# COM S 578X: Optimization for Machine Learning

Lecture Note 5: Optimality Conditions

Jia (Kevin) Liu

Assistant Professor
Department of Computer Science
Iowa State University, Ames, Iowa, USA

Fall 2019

# Recap Last Lecture

Given a minimization problem

> Minimize $\quad f(\mathbf{x})$
>
> subject to $\quad g_i(\mathbf{x}) \leq 0, \quad i = 1, \ldots, m \qquad \leftarrow u_i \geq 0$
>
> $\qquad\qquad h_j(\mathbf{x}) = 0, \quad j = 1, \ldots, p \qquad \leftarrow v_j$ unconstrained

We define the Lagrangian:

$$L(\mathbf{x}, \mathbf{u}, \mathbf{v}) = f(\mathbf{x}) + \sum_{i=1}^{m} u_i g_i(\mathbf{x}) + \sum_{j=1}^{p} v_j h_j(\mathbf{x})$$

and the Lagrangian dual function:

$$\Theta(\mathbf{u}, \mathbf{v}) = \min_{\mathbf{x}} L(\mathbf{x}, \mathbf{u}, \mathbf{v})$$

# Recap Last Lecture

The subsequent Lagrangian dual problem is:

$$\text{Maximize} \quad \Theta(\mathbf{u}, \mathbf{v})$$
$$\text{subject to} \quad \mathbf{u} \geq \mathbf{0}$$

Important properties:

- Dual problem is always convex (or $\Theta$ is always concave), even if the primal problem is non-convex
- The weak duality property always holds, i.e., the primal and dual optimal values $p^*$ and $d^*$ satisfy $p^* \geq d^*$
- Slater's condition: for convex primal, if $\exists \mathbf{x}$ such that

$$g_1(\mathbf{x}) < 0, \ldots, g_m(\mathbf{x}) < 0 \text{ and } h_1(\mathbf{x}) = 0, \ldots, h_p(\mathbf{x}) = 0.$$

then strong duality holds: $p^* = d^*$.

# Outline

Today:

- KKT conditions

- Geometric interpretation

- Relevant examples in machine learning and other areas

# Karush-Kuhn-Tucker Conditions

Given general problem

Minimize $\quad f(\mathbf{x})$

subject to $\quad g_i(\mathbf{x}) \leq 0, \quad i = 1, \ldots, m \qquad \leftarrow u_i \geq 0$

$\qquad\qquad\quad h_j(\mathbf{x}) = 0, \quad j = 1, \ldots, p \qquad \leftarrow v_j$ unconstrained

$$L(\underline{x}, \underline{u}, \underline{v}) = f(\underline{x}) + \underline{u}^T g(\underline{x}) + \underline{v}^T h(\underline{x})$$

The Karush-Kuhn-Tucker (KKT) conditions are:

*the grad of $L(x, \underline{u}, \underline{v})$ w.r.t. $\underline{x}$ is $\underline{0}$*

- Stationarity (ST): $\nabla_{\mathbf{x}} f(\mathbf{x}) + \sum_{i=1}^m u_i \nabla_{\mathbf{x}} g_i(\mathbf{x}) + \sum_{j=1}^p v_j \nabla_{\mathbf{x}} h_j(\mathbf{x}) = 0$
- Complementary slackness (CS): $u_i g_i(\mathbf{x}) = 0, \ \forall i$ *either $u_i = 0$ or $g_i(\underline{x}) = 0$*
- Primal feasibility (PF): $g_i(\mathbf{x}) \leq 0, \ h_j(\mathbf{x}) = 0, \ \forall i, j$
- Dual feasibility (DF): $u_i \geq 0, \ \forall i$

# KKT Necessity

$\mathbf{x}^*$ primal opt. $\left.\begin{array}{l}\end{array}\right\}$ strong duality $\Longrightarrow$ $(\mathbf{x}^*, \mathbf{u}^*, \mathbf{v}^*)$ are KKT.
$(\mathbf{u}^*, \mathbf{v}^*)$ dual opt.

> **Theorem 1**
>
> If $\mathbf{x}^*$ and $\mathbf{u}^*, \mathbf{v}^*$ be primal and dual *optimal* solutions w/ zero duality gap (e.g., implied by convexity and Slater's condition), then $(\mathbf{x}^*, \mathbf{u}^*, \mathbf{v}^*)$ satisfy KKT conditions.

WTS: (ST), (CS), (PF), (DF)

*Proof.* We have PF and DF for free from the assumption. Also, $\mathbf{x}^*$ and $(\mathbf{u}^*, \mathbf{v}^*)$ are primal & dual solutions with strong duality $\Rightarrow$

strong duality    def of dual fn.

$$f(\mathbf{x}^*) = \Theta(\mathbf{u}^*, \mathbf{v}^*) = \min_{\mathbf{x}} \left\{ f(\mathbf{x}) + \sum_{i=1}^{m} u_i^* g_i(\mathbf{x}) + \sum_{j=1}^{p} v_j^* h_j(\mathbf{x}) \right\}$$

DF PF    PF

def of "min" (=).

$\leq L$ evaluated at $\mathbf{x}$ any $\mathbf{x}$ (including $\mathbf{x}^*$)

$$\leq f(\mathbf{x}^*) + \sum_{i=1}^{m} \underbrace{u_i^* g_i(\mathbf{x}^*)}_{\leq 0} + \sum_{j=1}^{p} \underbrace{v_j^* h_j(\mathbf{x}^*)}_{=0} \leq f(\mathbf{x}^*)$$

$\geq 0 \quad \leq 0$     $L(\mathbf{x}, \mathbf{u}, \mathbf{v})$

$\leq 0 \quad (=0)$

That is, all these inequalities are equalities. Then:

- $\mathbf{x}^*$ minimizes $L(\mathbf{x}, \mathbf{u}^*, \mathbf{v}^*)$ over $\mathbf{x} \in \mathbb{R}^n$ (unconstrained) $\Rightarrow$ <u>Gradient of $L(\mathbf{x}, \mathbf{u}^*, \mathbf{v}^*)$ must be 0 at $x^*$</u>, i.e., the <span style="color:red">stationarity</span> condition. (ST)
- Since $u_i^* g_i(\mathbf{x}^*) \leq 0$ (PF & DF), we must have each $\underline{u_i^* g_i(\mathbf{x}^*) = 0}$, i.e., <span style="color:red">complementary slackness</span> condition. (CS) $\quad \square$

# KKT Sufficiency

If $(\underline{x}^*, \underline{u}^*, \underline{v}^*)$ is KKT $\Big\}$ $\Rightarrow$ $\Big\{$ $\underline{x}^*$ is primal opt.

primal is convex $\qquad\qquad$ $(\underline{u}^*, \underline{v}^*)$ is dual opt.

## Theorem 2

*If the primal problem is convex and $\mathbf{x}^*$ and $(\mathbf{u}^*, \mathbf{v}^*)$ satisfy KKT conditions, then $\mathbf{x}^*$ and $(\mathbf{u}^*, \mathbf{v}^*)$ are primal and dual optimal solutions, respectively.*

*Proof.* If $\mathbf{x}^*$ and $(\mathbf{u}^*, \mathbf{v}^*)$ satisfy KKT conditions, then

Lagrangian: $L = f(\underline{x}) + \underline{u}^\mathsf{T} g(\underline{x}) + \underline{v}^\mathsf{T} h(\underline{x})$. From (ST): $\nabla_{\underline{x}} L(\underline{x}^*, \underline{u}^*, \underline{v}^*) = 0$

$$\Theta(\mathbf{u}^*, \mathbf{v}^*) \overset{(a)}{=} f(\mathbf{x}^*) + \sum_{i=1}^{m} u_i^* g_i(\mathbf{x}^*) + \sum_{j=1}^{p} v_j^* h_j(\mathbf{x}^*)$$

(ST)

$\underbrace{\qquad}_{=0 \text{ (CS)}}$ $\underbrace{\qquad}_{=0 \text{ (PF)}}$

$\underline{x}^*$ is a minimizer of $L(\underline{x}, \underline{u}^*, \underline{v}^*)$.

$$\overset{(b)}{=} f(\mathbf{x}^*), \qquad \underbrace{\qquad\qquad\qquad}_{=0}$$

where $(a)$ follows from ST and $(b)$ follows from CS.

Therefore, the duality gap is zero. Note that $\mathbf{x}^*$ and $(\mathbf{u}^*, \mathbf{v}^*)$ are PF and DF. Hence, they are primal and dual optimal, respectively. $\qquad\square$

# In Summary

So putting things together…

## Theorem 3

*For a convex optimization problem with strong duality (e.g., implied by Slater's conditions or other constraints qualifications):*

$$\mathbf{x}^* \text{ and } (\mathbf{u}^*, \mathbf{v}^*) \text{ are primal and dual solutions}$$
$$\Longleftrightarrow \mathbf{x}^* \text{ and } (\mathbf{u}^*, \mathbf{v}^*) \text{ satisfy KKT conditions}$$

Warning: This statement is only true for convex optimization problems. For non-convex optimization problems, KKT conditions are neither necessary nor sufficient! (more on this shortly)

# Where Does This Name Come From?

Older books/papers referred to this as the KT (Kuhn-Tucker) conditions

- First appeared in a publication by Kuhn and Tucker in 1951
- Kuhn & Tucker shared the John von Neumann Theory Prize in 1980
- Later people realized that Karush had the same conditions in his unpublished master's thesis in 1939,



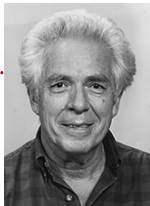William **K**arush

*Univ. of Chicago
Ph.D. advisor:
Magnus Hestenes.
(Conj. GD).*

*Lipschitz*

*Gauss.
'75*

Harold W. **K**uhn

*Princeton:
Friend with
John Nash
Hungarian Alg.
$O(n^4)$ $(O(n^3))$.*

Albert W. **T**ucker

*Princeton:
John Nash
Lloyd Shapley
(Shapley
value...
stochastic
games).
Marin Minsky*

- A Fun Read: R. W. Cottle, "William Karush and the KKT Theorem," Documenta Mathematica, 2012, pp. 255-269.

*RAND.
Friend with Richard
Bellman. → Assoc. Prof
at Chicago. → Cal State
University → Died
1997.*

# Other Optimality Conditions

- KKT conditions are a special case of the more general Fritz John Conditions:

$$\overset{\curvearrowleft}{u_0}\nabla f(\mathbf{x}^*) + \sum_{i=1}^{m} u_i \nabla g_i(\mathbf{x}^*) + \sum_{j=1}^{p} v_j \nabla h_j(\mathbf{x}^*) = \mathbf{0}$$
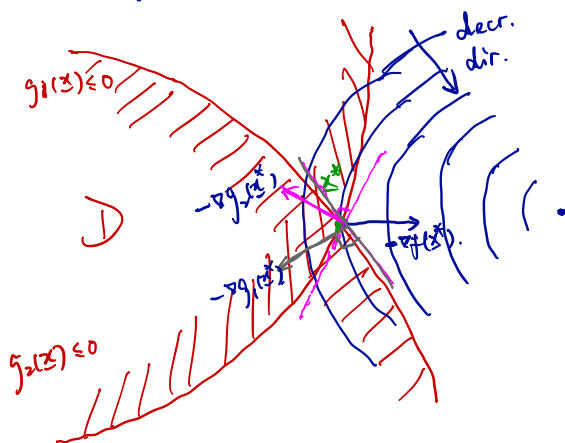
  where $u_0$ could be 0

- In turn, Fritz John conditions (hence KKT) belong to a wider class of the first-order necessary conditions (FONC), which allow for non-smooth functions using subderivatives

- Further, there are a whole class second-order necessary & sufficient conditiosn (SONC,SOSC) – also in "KKT style"

- For an excellent treatment on optimality conditions, see [BSS, Ch.4–Ch.6]

# Geometric Interpretation of KKT

Set of binding (active/tight) constraints: $\mathcal{I}(z^*) \triangleq \{i : g_i(z^*) = 0\}$.

(CS): $u_i^* \cdot g_i(z^*) = 0 \implies u_i^* \geq 0$

(ST): $\nabla f(z^*) + \sum_i u_i^* \nabla g_i(z^*) + \sum_j v_j^* \nabla h_j(z^*) = 0$



$g_1(z) \leq 0$

$g_2(z) \leq 0$

decr. dir.

$-\nabla g_1(z^*)$

$-\nabla g_2(z^*)$

$-\nabla f(z^*)$

$x^*$

physics interpretation

- $\nabla f(z^*)$: "pulling force".

- $\nabla g_i(z^*)$, $i \in \mathcal{I}(z^*)$.

  "repelling force".

$-\nabla g_1(z^*)$

$x^*$ $\longrightarrow$ $-\nabla f(z^*)$

$-g_2(z^*)$  sum $= 0$

# When is KKT neither sufficient nor necessary?

- (Not necc.): $\mathbf{x}^*$ is a (local) minimum $\not\Rightarrow$ $\mathbf{x}^*$ is a KKT point



$\mathbf{x}^*$ opt. but not KKT.

b/c : $\nabla f(\mathbf{x}^*) \neq u_1 \nabla g_1(\mathbf{x}^*) + u_2 \nabla g_2(\mathbf{x}^*)$.

$u_0 \quad \forall u_1, u_2 \geqslant 0$.

($\mathbf{x}^*$ is Fritz John pt. for $u_0 = 0$).

- (Not suff.): $\mathbf{x}^*$ is a KKT point $\not\Rightarrow$ $\mathbf{x}^*$ is a (local) minimum



obj: $\min \underline{c}^T \underline{x}$

$\underline{x}^*$ is KKT : $\exists \, u_1, u_2 \geqslant 0$

s.t. $-\underline{c} = u_1 \nabla g_1(\mathbf{x}^*) + u_2 \nabla g_2(\mathbf{x}^*)$

But $\underline{x}^*$ NOT opt.

# Example 1: Quadratic Problems with Equality Constraints

- Consider for $\mathbf{Q} \succeq 0$, the following quadratic programming problem is:

  Lagrangian: $\frac{1}{2}\mathbf{x}^\top \mathbf{Q}\mathbf{x} + \mathbf{c}^\top \mathbf{x} + \mathbf{u}^\top(\mathbf{A}\mathbf{x})$,

  $$\text{Minimize}_{\mathbf{x}} \quad \frac{1}{2}\mathbf{x}^\top \mathbf{Q}\mathbf{x} + \mathbf{c}^\top \mathbf{x}$$

  $$\text{subject to} \quad \mathbf{A}\mathbf{x} = \mathbf{0} \qquad \leftarrow \mathbf{u}$$

- A convex problem w/o inequality constraints. By KKT, $\mathbf{x}$ is primal optimal iff

  (ST): $\mathbf{Q}\mathbf{x} + \mathbf{c} + \mathbf{A}^\top \mathbf{u} = \mathbf{0}$

  (PF): $\mathbf{A}\mathbf{x} = \mathbf{0}$

  (DF) & (CS): Implied by (PF)

  $$\Rightarrow \begin{bmatrix} \mathbf{Q} & \mathbf{A}^\top \\ \mathbf{A} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{u} \end{bmatrix} = \begin{bmatrix} -\mathbf{c} \\ \mathbf{0} \end{bmatrix}$$

  for some dual variable $\mathbf{u}$. A linear equation system combines ST & PF (CS and DF vacuous)

- Often arises from using Newton's method to solved equality-constrained problems $\{\min_{\mathbf{x}} f(\mathbf{x}) | \mathbf{A}\mathbf{x} = \mathbf{b}\}$

  By Taylor's SO expansion: $f(\mathbf{x}) \approx f(\bar{\mathbf{x}}) + \nabla f(\bar{\mathbf{x}})^\top(\mathbf{x} - \bar{\mathbf{x}}) + \frac{1}{2}(\mathbf{x} - \bar{\mathbf{x}})^\top H(\bar{\mathbf{x}})(\mathbf{x} - \bar{\mathbf{x}}) + o(\|\mathbf{x} - \bar{\mathbf{x}}\|^2)$

  Note: $\mathbf{A}\bar{\mathbf{x}} = \mathbf{b}$   $\mathbf{A}\mathbf{x} = \mathbf{b}$   $\Rightarrow$   $\mathbf{A}(\mathbf{x} - \bar{\mathbf{x}}) = \mathbf{0}$   $\Rightarrow$   $\mathbf{A}\tilde{\mathbf{x}} = \mathbf{0}$
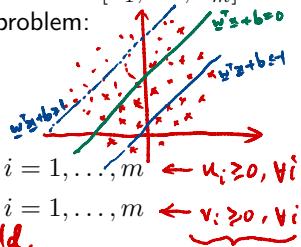
# Example 2: Support Vector Machine

Given labels $y \in \{-1, 1\}^n$, feature vectors $\mathbf{x}_1, \ldots, \mathbf{x}_m$. Let $\mathbf{X} \triangleq [\mathbf{x}_1, \ldots, \mathbf{x}_m]^\top$
Recall from Lecture 1 that the <span style="color:red">support vector machine</span> problem:

$$\underset{\mathbf{w}, b, \boldsymbol{\epsilon}}{\text{Minimize}} \quad \frac{1}{2}\|\mathbf{w}\|^2 + C\sum_{i=1}^m \epsilon_i$$

subject to $\quad$ (PF) $\begin{cases} y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \epsilon_i, & i = 1, \ldots, m \quad \leftarrow u_i \geq 0, \forall i \\ \epsilon_i \geq 0, & i = 1, \ldots, m \quad \leftarrow v_i \geq 0, \forall i \end{cases}$

*Slater's cond. hold.*
(PF)

Introducing dual variables $\mathbf{u}, \mathbf{v} \geq \mathbf{0}$ to obtain the KKT system:

Lagrangian: $\frac{1}{2}\|\mathbf{w}\|^2 + C\sum_{i=1}^m \epsilon_i + \sum_{i=1}^m u_i(1 - \epsilon_i - y_i(\mathbf{w}^\top \mathbf{x}_i + b)) - \sum_{i=1}^m v_i \epsilon_i$ $\quad$ *Quadratic in $\underline{w}$*

$\text{(ST)}: 0 = \sum_{i=1}^m u_i y_i, \quad \mathbf{w} = \sum_{i=1}^m u_i y_i \mathbf{x}_i, \quad \mathbf{u} = C\mathbf{1} - \mathbf{v}$ $\quad$ *Affine in $\underline{\epsilon}, b$*

$\text{(CS)}: v_i \epsilon_i = 0, \quad u_i(1 - \epsilon_i - y_i(\mathbf{x}_i^\top \mathbf{w} + b)) = 0, \quad i = 1, \ldots, m$
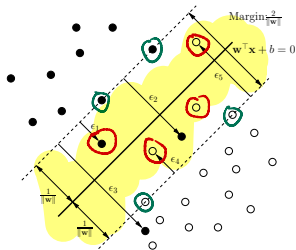
Take der.
w.r.t. : $\quad \underline{w}: \mathbf{w} - \sum_{i=1}^m u_i y_i \mathbf{x} = 0, \quad b: -\sum_{i=1}^m u_i y_i = 0, \quad \epsilon_i: C - u_i - v_i = 0, \forall i$

# Example 2: Support Vector Machine

$$(\text{ST}) \quad \underline{w} = \boxed{\underline{X} \, \text{Diag}\{y_1 \cdots y_m\}} \, \underline{u} = \underline{\tilde{X}} \, \underline{u}$$

Hence, at optimality, we have $\mathbf{w} = \sum_{i=1}^{m} u_i y_i \mathbf{x}_i$, and $u_i$ is nonzero only if $y_i(\mathbf{x}_i^\top \mathbf{w} + b) = 1 - \epsilon_i$. Such points are called the support points

- For support point $i$, if $\epsilon_i = 0$, then $\mathbf{x}_i$ lies on the edge of margin and $u_i \in (0, C]$ $\quad \epsilon_i = 0 \overset{(\text{CS})}{\Longrightarrow} v_i \geq 0 \overset{(\text{ST})}{\Longrightarrow} \underline{u} \leq C \underline{1}$

- For support point $i$, if $\epsilon_i \neq 0$, then $\mathbf{x}_i$ lies on wrong side of margin, and $u_i = C$ $\quad \epsilon_i \neq 0 \overset{(\text{CS})}{\Longrightarrow} v_i = 0 \overset{(\text{ST})}{\Longrightarrow} \underline{u} = C \underline{1}$



KKT conditions do not really give us a way to find solution here, but gives better understanding & useful in proofs

In fact, we can use this to screen away non-support points before performing optimization (lower-complexity)

# Constrained and Lagrange Forms

Often in ML and STATS, we'll switch back and forth between constrained form, where $t \in \mathbb{R}$ is a tuning parameter

$$\text{(C): } \min_{\mathbf{x}} f(\mathbf{x}) \quad \text{subject to} \quad g(\mathbf{x}) \leq t$$

*Special case: (t=0)*
*cannot find $\underline{x}$*
*s.t. $g(x) < 0$.*
*$g(x) = 0, \forall x$*
*set $u = \infty$*

and Lagrange form, where $u \geq 0$ is a tuning parameter

$$\text{(L): } \min_{\mathbf{x}} f(\mathbf{x}) + u \cdot g(\mathbf{x})$$

and claim these are equivalent. Is this true (assuming $f$ and $g$ convex)?

*WTS: $\underline{x}^* \in$ (C) w/ $t$ $\Rightarrow$ $\underline{x}^* \in$ (L) w/ $u$.*

*Proof.* (C) to (L): If Problem (C) is strictly feasible, then strong duality holds (why?), and there exists some $u \geq 0$ (dual solution) such that any solution $\mathbf{x}^*$ in (C) minimizes

*$\exists$ Dual var. $u$ s.t. $\underline{x}^*$ solves $f(x) + u(g(x) - t)$*

$$f(\mathbf{x}) + u \cdot (g(\mathbf{x}) - t). \quad = f(x) + u \cdot g(x) - ut$$

*const*

Clearly, $\mathbf{x}^*$ is also a solution in (L).

# Constrained and Lagrange Forms

(L) to (C): If $x^*$ is a solution in (L), then the KKT conditions for (C) are satisfied by taking $t = g(\mathbf{x}^*)$, so $\mathbf{x}^*$ is a solution in (C).

Putting things together: $(CS): u(g(\mathbf{x}^*) - t) = u(g(\mathbf{x}^*) - g(\mathbf{x}^*)) = 0.$

$$\bigcup_{u \geq 0} \left\{ \text{solutions in (L)} \right\} \quad \subseteq \quad \bigcup_{t} \left\{ \text{solutions in (C)} \right\}$$

$$\bigcup_{u \geq 0} \left\{ \text{solutions in (L)} \right\} \quad \supseteq \quad \bigcup_{\substack{t: \text{ (C) is strictly} \\ \text{feasible}}} \left\{ \text{solutions in (C)} \right\}$$

WTS: $x^* \in (L)$ w/ $u \Rightarrow x^* \in (C)$ w/ $t$, i.e., Given $u \geq 0$, WTF $\exists t$, s.t. KKT is satisfied

Try $t = g(x^*)$, check: (ST) $x^*$ is soln to (L) $\Rightarrow \nabla f(x^*) + u \nabla g(x^*) = 0$. (PF) & (DF)

I.e., nearly perfect equivalence. Note: If the only value of $t$ that leads to a feasible but not strictly feasible constraint set is $t = 0$, then we do get perfect equivalence

So, e.g., if $g \geq 0$ and (C) and (L) are feasible for all $t, u \geq 0$, then we do get perfect equivalence

$$\left. \begin{array}{l} g(x) \geq 0 \\ g(x) \leq t \end{array} \right\} \Rightarrow \begin{array}{l} 1^\circ \text{ when } t \neq 0, \text{ (C) is strictly feas.} \\ 2^\circ \text{ when } t = 0, \ g(x) = 0. \end{array} \right\} \begin{array}{l} \text{perfect equivalence.} \\ \text{If } g(x) \text{ is some norm.} \end{array}$$

# Next Class

Gradient Descent