# Back to Optimization with Constraints

$$\min \quad f(x)$$

$$\text{s.t} \quad g_i(x) \leq 0 \qquad i = 1 \cdots m$$

$$h_j(x) = 0 \qquad j = 1 \cdots \ell$$

*Temporarily absorb $h_j$'s as 2 ineq constraints*

*We need $-h_j$ & $h_j$ both convex & ∴ affine $h_j$*

from midsem Question

① Define: $I_{g_i}(x) = 0$ if $g_i(x) \leq 0$ &

$$= \infty \quad \text{o/w}$$

$$\min_x \quad f(x) + \sum_i I_{g_i}(x) \quad :- \text{Convex function But not diff.}$$

Solve by either analysing optimality conditions in terms of subgradients or employ subgradient descent...

② Write it equivalently as a cone program (yet to be analysed)...MIDSEM

③ Replace $\mathbb{I}_{g_i}(x)$ with a more "graceful" penalty function

$$\min_{x} f(x) - \sum_{i=1}^{m} \lambda_i \log(-g_i(x))$$

iteratively decrease $\lambda_i \geq 0$

④ Instead consider the Lagrangian fn

$$L(x,\lambda) = f(x) + \sum \lambda_i g_i(x)$$

We will briefly visit ① & then ④ & later ② & ③

⑤ Recall gradient descent & Newton:

$$x^{k+1} = \min_{x} f(x^k) + \nabla^T f(x^k)(x - x^k) + \frac{t}{2}(x - x^k)^T M (x - x^k)$$

$M = I$ for gradient desc & $\nabla^2 f(x^k)$ for Newton

"Proximal"/"Mirror descent"/Projection algos treat problem of finding $x^{k+1}$ as that of locating next iterate ==as close as possible to $x^k$==

In the sense of an approximation or in the sense of minimizing constraint violation etc

Suppose!
$$\min f(x)$$
$$\text{s.t} \quad g_i(x) \leq 0$$

$$\min f(x) + \eta \max_i g_i(x)$$

(we let $\eta$ iteratively tend to $\infty$)

You need to find the formulation of the constrained opt problem for which the subgradient can be discovered easily.

H/W   Eg Lasso: $\min_x \|Ax - y\|_2^2 \longrightarrow$ Regression loss/error

$$\|x\|_1 \leq \theta$$

$$\min f(x)$$
$$st \quad g_i(x) \le 0$$

**option 1. (0/1)**

$$I(x) \atop g_i \searrow 0 \quad \text{if } g_i(x) \le 0$$
$$\searrow \infty \quad o/w$$

if $g_i$ is convex, dmn $I_{g_i}$ is convex & $I_{g_i}(x)$ is a convex fn

$$\partial I_{g_i}(x) = \left\{ d \in \mathbb{R}^n \,\middle|\, I_{g_i}(y) \ge I_{g_i}(x) + d^T(y-x) \,\forall y \right\}$$

$\infty$ if $g_i(y) > 0$
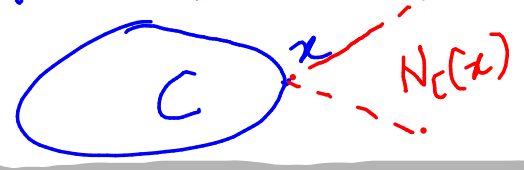so no issues

$0$ if $g_i(y) \le 0$

(if $g_i(x) \le 0$) $\left\{ d \in \mathbb{R}^n \,\middle|\, 0 \ge d^T(y-x) \,\forall y \text{ s.t. } g_i(y) \le 0 \right\}$

$$= \left\{ d \in \mathbb{R}^n \,\middle|\, d^T x \ge d^T y \,\forall y \text{ s.t. } g_i(y) \le 0 \right\}$$

Normal cone $N_c(x)$ to C at point $x$.
① If $x \in int(dmn\ g_i)$ then $N_c(x) = \{0\}$ i.e no nontrivial descent possible ② Otherwise

$$N_c(x) = \left\{ d \in \mathbb{R}^n \,\middle|\, d^T x \ge d^T y \,\forall y \in C \right\}$$



**option 2 (continuous)**

Let $C_i = \{ x \mid g_i(x) \le 0 \}$ are convex sets & let

$$dist(x, C_i) = \min \{ \|x - u\| : u \in C \}$$

If $C_i$ is closed, convex then

$\exists$ unique $u^* \in C$ that minimizes $\|x - u\|$. Let us call $u^* = P_{C_i}(x)$ so that

$$dist(x, C_i) = \|x - P_{C_i}(x)\|$$

We are interested in

$\hat{x}$ st $g_1(x) \le 0, \dots g_m(x) \le 0$

ie $\hat{x} \in C_1 \cap C_2 \dots \cap C_m$

Claim: (if $\hat{x}$ exists)

$$\min_{x \in \mathbb{R}^n} \ \max_{i=1\dots m} dist(x, C_i) = 0$$

call it $D(x)$

$D(\hat{x}) = 0$

$$\nabla dist(x, C_i) = \frac{x - P_{C_i}(x)}{\|x - P_{C_i}(x)\|}$$

if $D(x) = dist(x, C_i) \ne 0$ then
$$\frac{x - P_{C_i}(x)}{\|x - P_{C_i}(x)\|} \in \partial D(x)$$

For Lasso, it can be shown that for every $\theta$ there exists a $\lambda \geq 0$ s.t following two problems are equivalent:

$$\min_x \|Ax - y\|_2^2 + \lambda \|x\|_1$$

① ... say soln is $x^*$ & $\|x^*\|_1 = \beta$

$$\min_x \|Ax - y\|_2^2$$
$$\text{s.t } \|x\|_1 \leq \theta$$

② ... say solution is $\hat{x}$

Solution to ② with $\theta = \beta = x^*$ is also $x^*$!

Solution to ① with $\lambda$ as soln to $A^T(y - Ax) = \lambda g_{\hat{x}}$ is also $\hat{x}$!

$$g_{\hat{x}} \in \partial \|\hat{x}\|_1$$

**H/W:** Subgradient of $\|x\|_1 = f(x)$ $\qquad x \in \mathbb{R}^n$

$$f(x) = \|x\|_1 = \max_{i=1\ldots N} \left\{ f_1(x), f_2(x) \ldots f_i(x) \ldots f_N(x) \right\}$$

$$N = 2^n$$

$$S_1^T x \qquad S_2^T x \qquad\qquad S_N^T x$$

$$S_1 = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \\ \vdots \end{bmatrix} \qquad S_2 = \begin{bmatrix} -1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \qquad S_N = \begin{bmatrix} -1 \\ -1 \\ \vdots \\ -1 \end{bmatrix}$$

if No component of $x = 0$ then $\quad S = \begin{bmatrix} sgn(x_1) \\ sgn(x_2) \\ \cdot \\ sgn(x_n) \end{bmatrix}$

In general if $f(x) = S_1'^T x = S_2'^T x = \ldots = S_k'^T x$

$\left\{\begin{array}{l} +1 \quad \text{if } x_i > 0 \\ -1 \quad \text{if } x_i < 0 \end{array}\right.$

then $\partial f(x) = conv\left\{ S_1', S_2' \ldots S_k' \right\} \ldots (\partial f(x))_i =$

$\theta(-1) + (1-\theta)(1) \text{ if } x_i = 0$

$\theta \in [0,1]$

$\stackrel{ie}{=} \partial f(x) = \left\{ d \,\Big|\, \|d\|_\infty \leq 1, \; d^T x = \|x\|_1 \right\}$

**Claim:** If $\bar{a}_j = 2\sum_{i=1}^{n}[A_{ij}]^2$ & $\bar{b}_j = 2\sum_{i=1}^{n}A_{ij}(y_i - x_j^{T_n}A_j)$

**Then:** 
$$x_j^* = \begin{cases} (\bar{b}_j + \lambda)/\bar{a}_j & \text{if } \bar{b}_j < -\lambda \\ 0 & \text{if } \bar{b}_j \in [-\lambda, \lambda] \\ (\bar{b}_j - \lambda)/\bar{a}_j & \text{if } \bar{b}_j > \lambda \end{cases}$$

So to satisfy this, Lasso iterates on $x^k$ as follows: $x^{(0)} \to \bar{b}^{(0)} \to x^{(1)} \ldots \ldots - \text{until}$ convergence $\ldots$ We can understand through following simplification where $A = I$

---

**Eg:** $\min_{x} \frac{1}{2}\|y - x\|^2 + \underset{\text{Regularizer}}{\lambda\|x\|_1}$  $\left(\operatorname{argmin}_{x} \|y - x\|^2 + \lambda\|x\|_1 = x^*\right)$

I will suggest a soln by setting "some" $g_x = 0$   $\lambda \geq 0$

Higher $\lambda \Rightarrow$ more $x_i$'s are zeros ⊛

$$x_i^* = \begin{cases} -\lambda + y_i & \text{if } y_i > \lambda \\ 0 & \text{if } -\lambda \leq y_i \leq \lambda \\ \lambda + y_i & \text{if } y_i < -\lambda \end{cases}$$

lots of zeros esp if several $|y_i| < \lambda \ldots$ sparsity

Why should this be imp for minimization? 2 ways of answering

① $g_x = \frac{1}{2}\nabla(\|y - x\|^2) + \lambda\partial\|x\|_1$
$= (x - y) + \lambda\begin{bmatrix} \text{sign}(x_1) \\ \vdots \\ \text{sign}(x_n) \end{bmatrix}$

② $\min_{x_i} \frac{1}{2}(y_i - x_i)^2 + \lambda|x_i|$

for each $i$   $g_{x_i} = (x_i - y_i) + \lambda \text{sign}(x_i)$

In either case ① or ②, setting $g_x = 0$ or $g_{x_i} = 0$ for each $i$, & checking that ⊛ satisfies this equation,

---

## Another example

Maximum eigenvalue of a symmetric matrix

$$f(x) = \lambda_{max}(A(x)) \quad \because \quad A(x) = A_0 + x_1 A_1 + \ldots + x_n A_n$$
$$\& \quad A_i \in S^m$$

$$f(x) = \lambda_{max}(A(x)) = \sup \quad y^T A(x) y$$

Index set $I$ over $\|y\|_2 = 1$
fns

each function is affine in $x$ for fixed $y$ & has gradient $\nabla f_y(x) = (y^T A_1 y, \ldots y^T A_n y)$

http://en.wikipedia.org/wiki/Rayleigh_quotient

Active fns $y^T A(x) y$ are the ones for which $y$ is (normalised) eigenvector for max eigenvalue $\lambda_{max}$ of $A(x)$

$$\therefore g_x = (y^T A_1 y, \ldots y^T A_n y)$$

④

$$L(x, \lambda, \mu) = f(x) + \sum_i \lambda_i g_i(x) + \sum_j \mu_j h_j(x)$$

for

$$\min f(x)$$
$$\text{s.t} \quad g_i(x) \leq 0 \quad i = 1 \ldots m$$
$$h_j(x) = 0 \quad j = 1 \ldots l$$

$\min L(x, \lambda, \mu)$ --- you should ideally have $\lambda_i \geq 0$
to penalize $g_i(x) > 0$

$$\min_x f(x) \geq \min_x f(x) + \sum_i \lambda_i g_i(x) + \sum_j \mu_j h_j(x) \geq \min_{x, \lambda_i \geq 0} L(x, \lambda, \mu)$$
$$g_i(x) \leq 0 \qquad g_i(x) \leq 0, \ \lambda_i \geq 0$$
$$h_j(x) = 0 \qquad h_j(x) = 0$$

$$\min_{x} f(x) \geq \boxed{\max_{\substack{\lambda \geq 0 \\ \mu_j}} \min_{x} L(x, \lambda, \mu)}$$

$$\text{s.t } g_i(x) \leq 0$$

$$h_j(x) = 0$$

Pushes up the lower bound from previous inequality.

Importance: ① maximizing $L^*(\lambda, \mu)$ subject to just $\lambda_i \geq 0$ could be easier & provide a lower bound for original objective, provided $L^*(\lambda, \mu)$ has a manageable form!

② $L^*(\lambda, \mu)$ is min over affine fns $L(x, \lambda, \mu)$ indexed by $x \ldots$ ∴ $L^*(\lambda, \mu)$ is concave fn of $\lambda$ & $\mu$

Max over concave + nice ∴ properties that min over convex has

③ $L^*(\lambda, \mu)$ is concave for all $f$, $g_i$ & $h_j$ choices. HOWEVER if $f$ is convex, $g_i$'s are convex & $h_j$'s affine, most often the lower bound $\left(\max\limits_{\substack{\lambda_i \geq 0 \\ \mu_j}} L^*(\lambda, \mu)\right)$ turns out to be the exact solution..

Duality gap: $f(x^*) - L^*(\lambda^*, \mu^*)$

[Eg: Support Vector m/cs have dual easier to solve very often & dual gives exact solution as primal]

Soln: $x^a$

PRIMAL: $\min\limits_{x} f(x)$
st $g_i(x) \leq 0$
$h_j(x) = 0$

$\geq$

(equality under convexity++)

DUAL: $\max\limits_{\substack{\lambda_i \geq 0 \\ \mu_j}} L^*(\lambda, \mu)$

Soln: $\lambda_i^a$, $\mu_j^a$

④ Convergence criterion: Say Primal $\left(\boxed{f(x^*)}\right) = $ Dual $\left(\boxed{L^*(\lambda^a, \mu^a)}\right)$

$f(x^k) - L^*(\lambda^k, \mu^k)$ can be used as measure of distance from optimal soln where duality gap $= 0$

Recall: $f(x^*) \geq L^*(\lambda, \mu)$ $\forall \lambda, \mu$

$$\min_x f(x) \leq \max_{\substack{\lambda_i \geq 0 \\ \mu \in R}} L^{*}(\lambda, \mu)$$

$$s.t \quad g_i(x) \leq 0$$

$$h_j(x) = 0$$

**Q1:** Did we require $f$, $g_i$'s & $h_j$'s to be convex or affine? Ans: No

**Q2:** Is $L^{*}$ concave irrespective of $f$, $g_i$'s & $h_j$'s? Note: $L(x, \lambda, \mu)$ is affine in $\lambda, \mu$

$$L^{*} = \min_x \underbrace{L(x, \lambda, \mu)}_{L_x(\lambda, \mu)}$$



min of affine fns is concave

Next we provide insight into equality between primal & dual solutions through min-max theorem & saddle point theorem

**Claim 1:** Min-max / inf·sup inequality

$$\sup_y \inf_x f(x,y) \leq \inf_x \sup_y f(x,y)$$

**Proof:** $\forall x \left[ f(x,y) \leq \sup_y f(x,y) \right] \Rightarrow \left[ \inf_x f(x,y) \leq \inf_x \sup_y f(x,y) \right]$

$$\forall y$$
$$\Downarrow$$
$$\sup_y \inf_x f(x,y) \leq \inf_x \sup_y f(x,y)$$

**Claim 2**

.. A saddle point of a function $f : \mathcal{X} \times \mathcal{Y} \to \Re \cup \{\pm \inf\}$ is a pair $(\overline{x}, \overline{y}) \in \mathcal{X} \times \mathcal{Y}$ satisfying

$$\sup_y f(\overline{x}, y) \leq f(\overline{x}, \overline{y}) \leq \inf_x f(x, \overline{y})$$

Show that if $f(x,y)$ has a saddle point $(\overline{x}, \overline{y})$ then

$$\sup_y \inf_x f(x,y) = \inf_x \sup_y f(x,y)$$

**Proof:** By the min-max / inf-sup inequality

$$\sup_y \inf_x f(x,y) \leq \inf_x \sup_y f(x,y) \to \textcircled{1}$$

Now, if $f$ has a saddle point $(\overline{x}, \overline{y})$ then

By defn of saddle pt

$$\inf_x \sup_y f(x,y) \leq \sup_y f(\overline{x}, y) \leq f(\overline{x}, \overline{y}) \leq \inf_x f(x, \overline{y})$$
$$\leq \sup_y \inf_x f(x,Y)$$

1

Thus

$$\inf_x \sup_y f(x,y) \le \sup_y \inf_x f(x,y) \rightarrow ②$$

By ① & ②, we have min-max equality!

$$\sup_y \inf_x f(x,y) = \inf_x \sup_y f(x,y) = f(\bar{x}, \bar{y})$$

Illustration of saddle point at $(0,0)$ for $f(x_1, x_2) = x_1^2 - x_2^2$
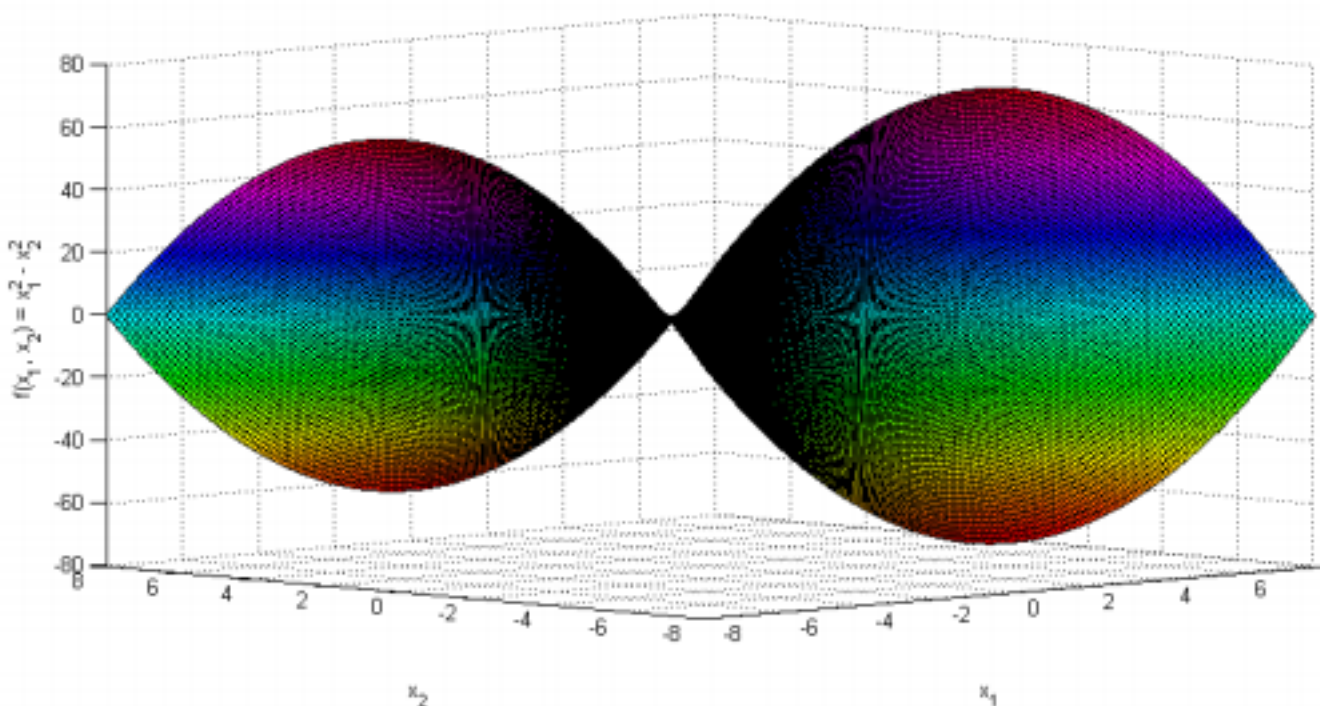on pages 242 & 243 of http://www.cse.iitb.ac.in/~cs709/notes/BasicsOfConvexOptimization.pdf



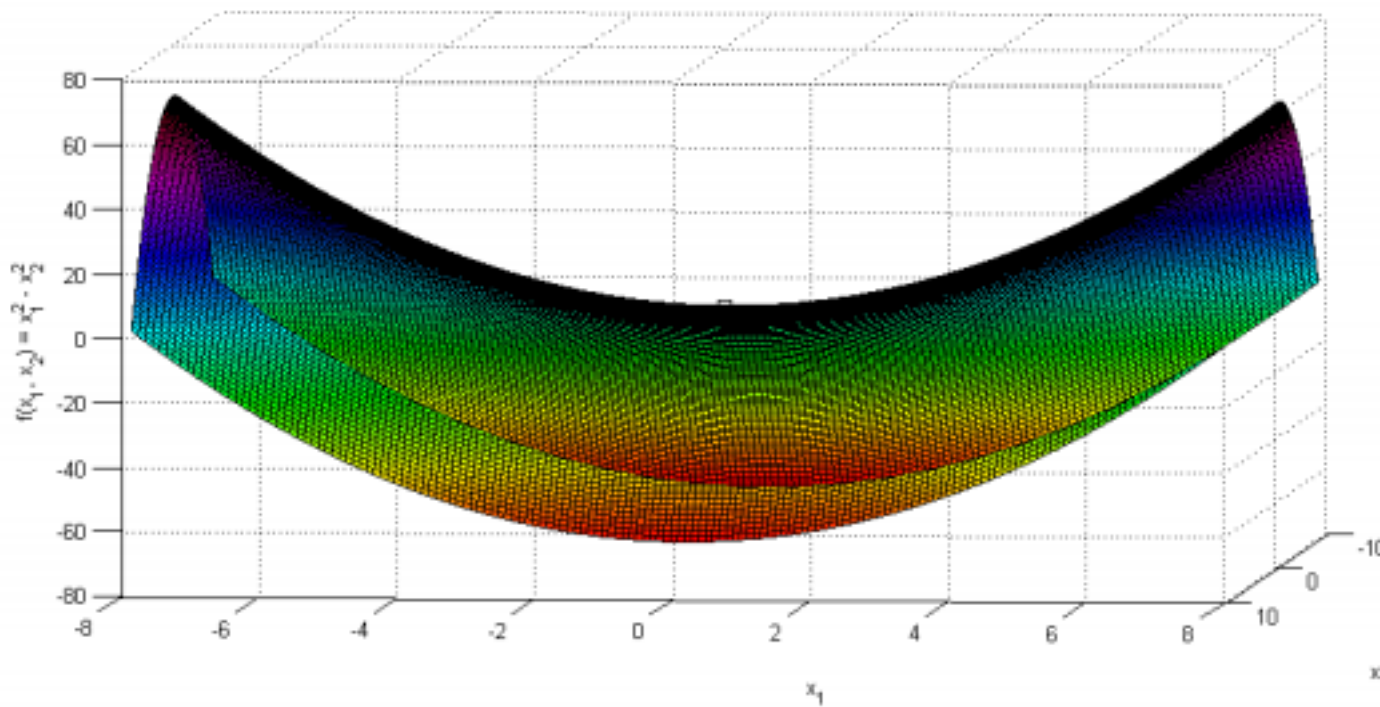Figure 4.19: The hyperbolic paraboloid $f(x_1, x_2) = x_1^2 - x_2^2$, which h

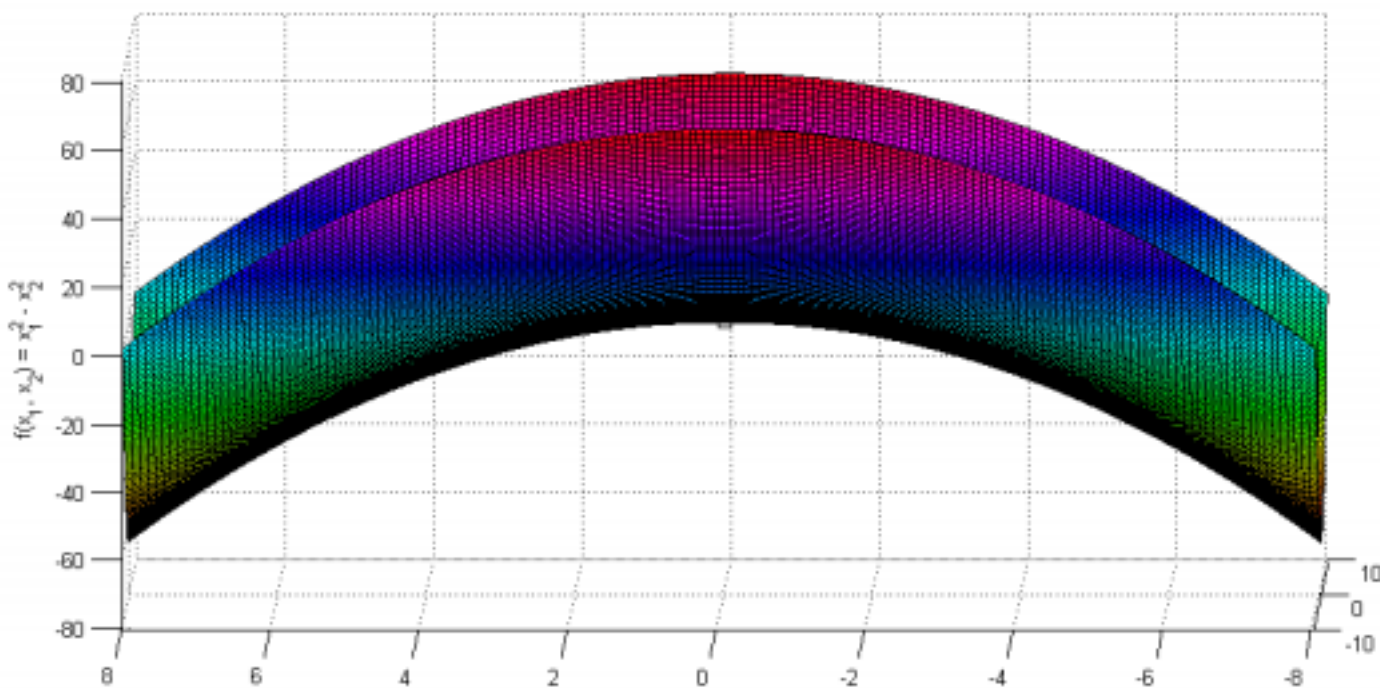Figure 4.20: The hyperbolic paraboloid $f(x_1, x_2) = x_1^2 - x_2^2$, when vie the $x_1$ axis is concave up.



Figure 4.21: The hyperbolic paraboloid $f(x_1, x_2) = x_1^2 - x_2^2$, when vie the $x_2$ axis is concave down.

$$\begin{cases} \min \quad f(x) \\ s.t \quad g_i(x) \le 0 \quad i=1\ldots m \\ \quad h_j(x)=0 \quad j=1\ldots k \end{cases}$$

We will generalize the inequalities & equalities

$$\min_x f(x)$$
$$s.t \quad g_i(x) \le 0$$
$$h_j(x)=0$$

$$\ge \min_x \max_{\lambda, \mu} f(x) + \sum_{i=1}^{m} \lambda_i g_i(x) + \sum_{j=1}^{k} \mu_j h_j(x)$$
$$s.t \quad g_i(x) \le 0$$
$$h_j(x)=0$$
$$\lambda_i \ge 0 \quad \mu_j \in \mathbb{R}$$

$L(x, \lambda, \mu)$

By min-max thm: →

$\ge \min_x \max_{\lambda, \mu} L(x, \lambda, \mu)$
$\lambda_i \ge 0$
$\mu_j \in \mathbb{R}$

under strong duality
$\lambda_i^{*} g_i(x^{*}) = 0$
$\forall i$
$\lambda_i^{*} \ge 0$
$\mu_j^{*} h_j(x^{*}) = 0$
$\forall j$

By saddle pt thm: If $L(x, \lambda, \mu)$ has a saddle pt $(\bar{x}, \bar{\lambda}, \bar{\mu})$ then equality holds

$\ge \max_{\lambda, \mu} \min_x L(x, \lambda, \mu)$
$\lambda_i \ge 0$
$\mu_j \in \mathbb{R}$

General weak duality result

$L^{*}(\lambda, \mu)$ or lagrange dual fn.

$= \max_{\lambda, \mu, \lambda \ge 0} L^{*}(\lambda, \mu)$

Dual opt problem

# Karush Kuhn Tucker conditions (KKT conditions)

① These conditions can be used to determine a saddle point $(\bar{x}, \bar{\lambda}, \bar{\mu})$ for $L(x, \lambda, \mu)$

② We will first graphically motivate the KKT conditions

③ Next we will prove their necessity & sufficiency for optimality under convexity