**Definition 12.3.**

Given a feasible point $x$ and the active constraint set $\mathcal{A}(x)$ of Definition 12.1, the set of linearized feasible directions $\mathcal{F}(x)$ is

$$\mathcal{F}(x) = \left\{ d \;\middle|\; \begin{array}{ll} d^T \nabla c_i(x) = 0, & \text{for all } i \in \mathcal{E}, \\ d^T \nabla c_i(x) \geq 0, & \text{for all } i \in \mathcal{A}(x) \cap \mathcal{I} \end{array} \right\}.$$

Note : $\mathcal{F}(x)$ is a cone. $[\, C \subseteq \mathbb{R}^n$ is a cone if $x \in C, \lambda \geq 0 \Rightarrow \lambda x \in C.]$

From the three examples, we should expect that in general :

$x^*$ solves
$\min\limits_{x} f(x)$

s.t. $c_i(x) = 0, i \in \mathcal{E}$
$\quad c_i(x) \geq 0 \; i \in \mathcal{I}$

ideas from local linear approximation

$\Rightarrow \nexists \, d \in \mathbb{R}^n$ s.t. $d \in \mathcal{F}(x^*)$ and $\nabla f(x^*)^T d < 0$

and the remaining work is to convert this condition to a more convenient condition (one that involves Lagrange multipliers)

This step should only involves linear algebra!

We shall do exactly this to get our first major result of this course.

But notice an annoying technicality :

In our first example, if we change the constraint $\overbrace{x_1^2 + x_2^2 - 2 = 0}^{c_1(x)}$ to the equivalent :

$$\underbrace{(x_1^2 + x_2^2 - 2)^2}_{\text{new } c_1(x)} = 0, \quad \leftarrow \text{represent the same circle}$$

then

$$\nabla (\text{new } c_1)(x) = 2(x_1^2 + x_2^2 - 2) \begin{bmatrix} 2x_1 \\ 2x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \forall \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \text{ on the circle.}$$
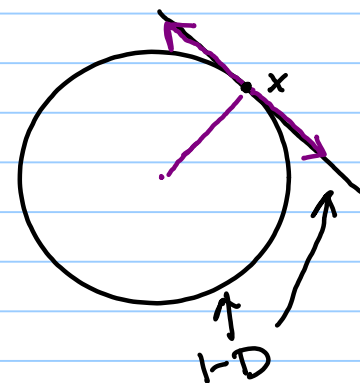
with the original $c_1$ :

$$\mathcal{F}(x) = \{d : 2x^T d = 0\} = \text{the tangent line}$$
$$= \text{all vectors} \perp x \qquad \text{of the constraint}$$
$$\text{set at } x.$$

with the new $c_1$ :

$$\mathcal{F}(x) = \{d : \begin{bmatrix} 0 \\ 0 \end{bmatrix}^T d = 0\} = \mathbb{R}^2 \leftarrow 2\text{-}D$$

1-D

Summary: $\min x_1 + x_2 \overset{= f(x)}{}$    has $x^* = \begin{bmatrix} -1 \\ -1 \end{bmatrix}$ as its solution, but the
      st. $\underbrace{(x_1^2 + x_2^2 - 2)^2}_{c_1(x)} = 0$    condition $\cancel{\circledast}$ does not hold!

We have actually seen this problem earlier:

In general, if $c(x) = 0$, we expect $\{y : c(y) = c(x)\} = \vec{c}^{-1}(c(x))$ to be
a hypersurface near $x$, $\nabla c(x)$ is orthogonal to the hypersurface,
and
$$\{d : \nabla c(x)^T d = 0\} = \text{the tangent plane of } \vec{c}^{-1}(c(x)) \text{ at } x.$$

$\uparrow$
$(n-1)$-dimensional when $\nabla c(x) \neq \vec{0}$

But this picture can totally fall apart if $\nabla c(x) = \vec{0}$! (recall my counter-
examples.)

So maybe everything is fine (for the expected optimality theorem) if we

impose      $\nabla c_i(x^*) \neq \vec{0} \;\; \forall i$ ?
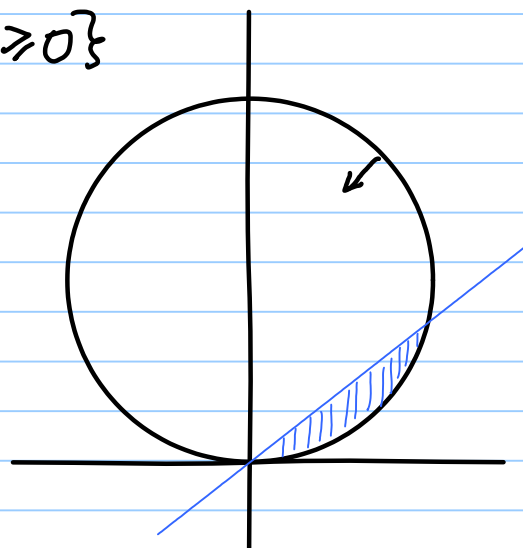
It turns out that we will still have a problem:

Consider the constraints $\quad c_1(x) = 1 - x_1^2 - (x_2-1)^2 \geq 0$

$\qquad\qquad\qquad\qquad\quad c_2(x) = -x_2 + mx_1 \qquad \geq 0 \quad , \; (m \in \mathbb{R})$

$\begin{bmatrix} 0 \\ 2 \end{bmatrix}$ $\qquad\qquad\qquad\qquad$ $\begin{bmatrix} m \\ -1 \end{bmatrix}$

$\mathcal{F}(\begin{bmatrix} 0 \\ 0 \end{bmatrix}) = \{ d \in \mathbb{R}^2 : \nabla c_1(\vec{0})^\top d \geq 0, \; \nabla c_2(\vec{0})^\top d \geq 0 \}$

$= $  $\cap$ 

$= $ 

good first order approximation near $\begin{bmatrix} 0 \\ 0 \end{bmatrix}$



what if $m = 0$?

$\{ c_1(x) \geq 0, \; c_2(x) \geq 0 \} = \{ \begin{bmatrix} 0 \\ 0 \end{bmatrix} \} \qquad \not\cong \qquad \mathcal{F}(\begin{bmatrix} 0 \\ 0 \end{bmatrix}) = \{ \begin{bmatrix} d_1 \\ 0 \end{bmatrix} : d_1 \in \mathbb{R} \}$

$\qquad\qquad\qquad\uparrow$ $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\uparrow$

$\qquad\qquad$ 0-dimensional $\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ 1-dimensional

If we consider $\min x_1 + x_2$ s.t. $\begin{cases} 1 - x_1^2 - (x_2-1)^2 \geq 0 \\ \qquad -x_2 \geq 0 \end{cases}$

The minimizer is at $\begin{bmatrix} 0 \\ 0 \end{bmatrix}$, as $\Omega = \left\{ \begin{bmatrix} 0 \\ 0 \end{bmatrix} \right\}$ ← the feasible region

The hoped-for necessary condition will not hold:

$$\underset{\nabla f(\begin{bmatrix} 0 \\ 0 \end{bmatrix})}{\underbrace{\begin{bmatrix} 1 \\ 1 \end{bmatrix}^T}} \underset{F(\begin{bmatrix} 0 \\ 0 \end{bmatrix})}{\underbrace{\begin{bmatrix} d_1 \\ 0 \end{bmatrix}}} < 0 \quad \forall d_1 < 0$$

In this case, neither $\nabla c_1(\vec{0})$ nor $\nabla c_2(\vec{0})$ is $\vec{0}$, but the two vectors are parallel. $\underset{\begin{bmatrix} 0 \\ 2 \end{bmatrix}}{} \qquad \underset{\begin{bmatrix} 0 \\ -1 \end{bmatrix}}{}$

**Definition 12.4** (LICQ).

    *Given the point $x$ and the active set $\mathcal{A}(x)$ defined in Definition 12.1, we say that the* linear independence constraint qualification (LICQ) *holds if the set of active constraint gradients* $\{\nabla c_i(x), i \in \mathcal{A}(x)\}$ *is* <u>*linearly independent.*</u>

**Theorem 12.1** (First-Order Necessary Conditions).

Suppose that $x^*$ is a local solution of (12.1), that the functions $f$ and $c_i$ in (12.1) are continuously differentiable, and that the LICQ holds at $x^*$. Then there is a Lagrange multiplier vector $\lambda^*$, with components $\lambda_i^*$, $i \in \mathcal{E} \cup \mathcal{I}$, such that the following conditions are satisfied at $(x^*, \lambda^*)$

$$\nabla_x \mathcal{L}(x^*, \lambda^*) = 0, \tag{12.34a}$$

$$c_i(x^*) = 0, \quad \text{for all } i \in \mathcal{E}, \tag{12.34b}$$

$$c_i(x^*) \geq 0, \quad \text{for all } i \in \mathcal{I}, \tag{12.34c}$$

$$\lambda_i^* \geq 0, \quad \text{for all } i \in \mathcal{I}, \tag{12.34d}$$

$$\lambda_i^* c_i(x^*) = 0, \quad \text{for all } i \in \mathcal{E} \cup \mathcal{I}. \tag{12.34e}$$

$\mathcal{L}(x, \lambda)$

$:= f(x) - \sum_{i \in \mathcal{E} \cup \mathcal{I}} \lambda_i c_i(x)$

if constraint $i$ is inactive (i.e. $c_i(x) > 0$), then the corresponding $\lambda_i = 0$

Example

$$\min_{x} \ (x_1 - 3/2)^2 + (x_2 - 1/2)^4 \quad \text{s.t.} \quad \begin{aligned} c_1 &= \\ c_2 &= \\ c_3 &= \\ c_4 &= \end{aligned} \begin{bmatrix} 1 - x_1 - x_2 \\ 1 - x_1 + x_2 \\ 1 + x_1 - x_2 \\ 1 + x_1 + x_2 \end{bmatrix} \geq 0 \qquad \mathcal{E} = \phi, \ \mathcal{I} = \{1, 2, 3, 4\}$$

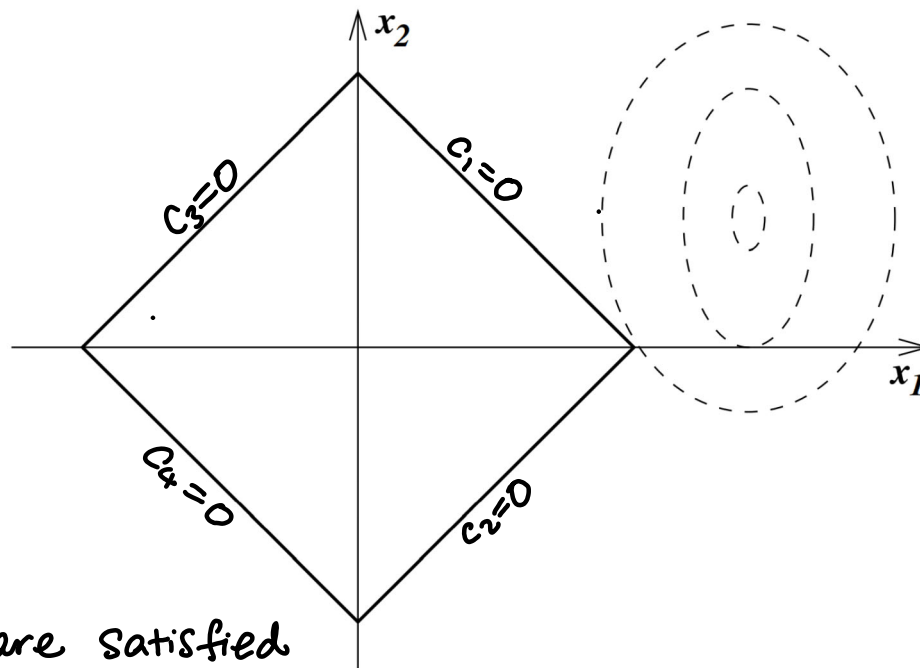Solution at $x^* = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$.

$\mathcal{A}(x^*) = \{1, 2\}$

$\nabla c_1(x^*) = \begin{bmatrix} -1 \\ -1 \end{bmatrix}$

$\nabla c_2(x^*) = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$

$\nabla f(x^*) = \begin{bmatrix} -1 \\ -1/2 \end{bmatrix} = \frac{3}{4} \begin{bmatrix} -1 \\ -1 \end{bmatrix} + \frac{1}{4} \begin{bmatrix} -1 \\ 1 \end{bmatrix}$

$\lambda^* = \begin{bmatrix} 3/4 \\ 1/4 \\ 0 \\ 0 \end{bmatrix}$.

KKT conditions are satisfied at $x^*$.

Strategy for proving KKT:

(I) we have explained that the hoped-for result

$$x^* \text{ solves } \min_x f(x) \text{ s.t. } c_i(x)=0, i\in\mathcal{E} \quad c_i(x)\geq 0 \; i\in\mathcal{I} \implies \nexists d\in\mathbb{R}^n \text{ st. } d\in\mathcal{F}(x^*) \text{ and } \nabla f(x^*)^T d < 0$$

$$\iff \nabla f(x^*)^T d \geq 0, \text{ for all } d\in\mathcal{F}(x^*)$$

does not hold without any assumption on $\nabla c_i(x)$, $i\in A(x^*)$.

Notice that the solution should depend only on $f$ and the feasible set $\Omega$ set itself, but not the algebraic specification of $\Omega$.

Yet, as we showed, the cone $\mathcal{F}(x^*)$ does depend on the algebraic specification of $\Omega$.

To correct this problem, we show that there is a more geometrically defined cone, called the tangent cone and denoted by $T_\Omega(x^*)$, so that

$$x^* \text{ is a solution} \implies \nabla f(x^*)^T d \geq 0, \text{ for all } d\in \overset{T_\Omega(x^*)}{\cancel{\mathcal{F}(x^*)}}$$

(II)  We show that under the LICQ assumption, $T_{lin}(x^*) = \mathcal{F}(x^*)$.

  This essentially follows from the implicit function theorem.

(III)  We show that: $\nabla f(x^*)^T d \geqslant 0$, for all $d \in \mathcal{F}(x^*)$
                      is equivalent to the KKT conditions.

  This essentially follows from Farkas' lemma.


Spirit :   • a minimizer $\Rightarrow$ a local minimizer

           • Steps (I) and (II) convert the local minimizer problem, using
             _local linear approximations_, into a linear algebra problem.

           • Step (III) is pure linear algebra.

## Definition 12.2.

*The vector $d$ is said to be a* tangent *(or tangent vector) to $\Omega$ at a point $x$ if there are a feasible sequence $\{z_k\}$ approaching $x$ and a sequence of <u>positive</u> scalars $\{t_k\}$ with $t_k \to 0$ such that*

$$\lim_{k \to \infty} \frac{z_k - x}{t_k} = d. \tag{12.29}$$

*The set of all tangents to $\Omega$ at $x^*$ is called the* tangent cone *and is denoted by $T_\Omega(x^*)$.*

$T_\Omega(x^*)$ depends only on the geometry of $\Omega$, not the algebraic specification of $\Omega$.

Recall:

## Definition 12.3.

*Given a feasible point $x$ and the active constraint set $A(x)$ of Definition 12.1, the set of* linearized feasible directions *$\mathcal{F}(x)$ is*

$$\mathcal{F}(x) = \left\{ d \;\middle|\; \begin{array}{ll} d^T \nabla c_i(x) = 0, & \text{for all } i \in \mathcal{E}, \\ d^T \nabla c_i(x) \geq 0, & \text{for all } i \in A(x) \cap \mathcal{I} \end{array} \right\}$$

depends on the algebraic specification of $\Omega$

The glory details.

(I) should be obvious intuitively, and is actually easy to prove.

**Theorem 12.3.**

*If $x^*$ is a local solution of (12.1), then we have*

means $f$ doesn't ↑ in direction $d$

means $d$ points towards $\llcorner\Omega$

$$\nabla f(x^*)^T d \geq 0, \quad \text{for all } d \in T_\Omega(x^*).$$

Proof: Assume the contrary that $\exists d \in T_\Omega(x^*)$ st. $\nabla f(x^*)^T d < 0$.

Let $\{z_k\}$ and $\{t_k\}$ be the sequences satisfying Definition 12.2 for this $d$.

Then:

$$f(z_k) = f(x^*) + \nabla f(x^*)^T \underbrace{(z_k - x^*)}_{= t_k d + o(t_k)} + \underbrace{o(\|z_k - x^*\|)}_{= o(t_k)}$$

$$= f(x^*) + \underbrace{\nabla f(x^*)^T d}_{<0} + o(t_k)$$

So, $f(z_k) < f(x^*) + 0.99 \nabla f(x^*)^T d$    for large enough $k$.

     ($z_k \to \Omega$)

This means $x^*$ cannot be a local solution. Q.E.D.

Of course, the converse of this result is not true.
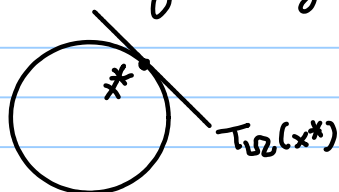
Counterexamples we have seen before :

(No constraint :) $\quad$ min $-\sum_i z_i^2$

$$\Omega = \mathbb{R}^n, \quad x^* = \vec{0}, \quad \nabla f(x^*) = \vec{0}$$

$$T_{L\Omega}(x^*) = \mathbb{R}^n \quad \nabla f(x^*)^T d = 0$$
$$\forall d.$$

But $x^*$ is not a *local* minimizer. ( It is a maximizer.)

(1 equality constraint :) $\quad$ min $x_1 + x_2 \quad$ s.t. $\quad x_1^2 + x_2^2 = 2$



$$x^* = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad T_{L\Omega}(x^*) = \{ d : \begin{bmatrix} 1 \\ 1 \end{bmatrix}^T d = 0 \}$$

maximizer! $\quad \nabla f(x^*) = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad \nabla f(x^*)^T d = 0 \quad \forall d \in T_{L\Omega}(x^*)$

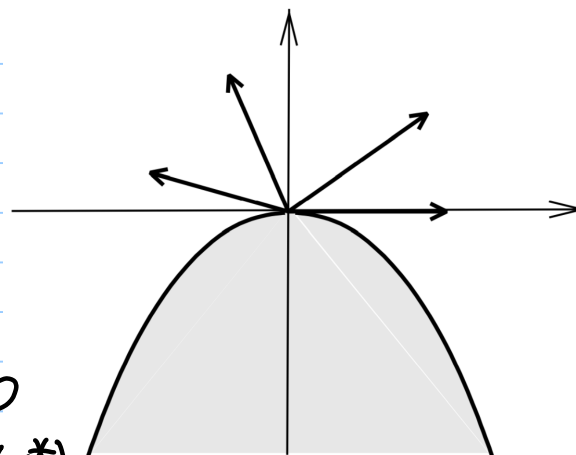A counterexample with one inequality constraint (less obvious) :

$$\min_{x \in \mathbb{R}^2} x_2 \quad \text{s.t.} \quad x_2 \geq -x_1^2$$

$$x^* = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \quad T_{L\Omega}(x^*) = \left\{ \begin{bmatrix} d_1 \\ d_2 \end{bmatrix} : d_2 \geq 0 \right\}$$

(why?)

$$\nabla f(x^*) = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \quad \nabla f(x^*)^T d = d_2 \geq 0$$
$$\forall d \in T_{L\Omega}(x^*)$$

But $x^*$ is not a local minimizer.

Step II

**Lemma 12.2.**

Let $x^*$ be a feasible point. The following two statements are true.

(i) $T_\Omega(x^*) \subset \mathcal{F}(x^*)$.

(ii) If the LICQ condition is satisfied at $x^*$, then $\mathcal{F}(x^*) = T_\Omega(x^*)$.

*Recall:*

**Theorem A.2** (Implicit Function Theorem).

Let $h : \mathbb{R}^n \times \mathbb{R}^m \to \mathbb{R}^n$ be a function such that

**(i)** $h(z^*, 0) = 0$ for some $z^* \in \mathbb{R}^n$,

**(ii)** the function $h(\cdot, \cdot)$ is continuously differentiable in some neighborhood of $(z^*, 0)$, and

**(iii)** $\nabla_z h(z, t)$ is nonsingular at the point $(z, t) = (z^*, 0)$.

Then there exist open sets $\mathcal{N}_z \subset \mathbb{R}^n$ and $\mathcal{N}_t \subset \mathbb{R}^m$ containing $z^*$ and $0$, respectively, and a continuous function $z : \mathcal{N}_t \to \mathcal{N}_z$ such that $z^* = z(0)$ and $h(z(t), t) = 0$ for all $t \in \mathcal{N}_t$. Further, $z(t)$ is uniquely defined. Finally, if $h$ is $p$ times continuously differentiable with respect to both its arguments for some $p > 0$, then $z(t)$ is also $p$ times continuously differentiable with respect to $t$, and we have

$$\nabla z(t) = -\nabla_t h(z(t), t) [\nabla_z h(z(t), t)]^{-1}$$

for all $t \in \mathcal{N}_t$.

Proof of (i): Let $d \in T_{\Omega}(x^*)$. By definition, $\exists \ z_k \in \Omega$, $t_k > 0$ s.t.

$$z_k \to x^*, \ t_k \to 0 \quad \text{and} \quad \frac{z_k - x^*}{t_k} \to d.$$

$$\Updownarrow$$

$$\frac{z_k - x^* - t_k d}{t_k} \to 0 \iff z_k = x^* + t_k d + o(t_k).$$

If $i \in A(x^*) \cap \mathcal{E}$,
then

$$0 = \frac{1}{t_k} c_i(z_k) = \frac{1}{t_k}\Big[ \underbrace{c_i(x^*)}_{=0} + \nabla c_i(x^*)^\top \underbrace{(z_k - x^*)}_{t_k d + o(t_k)} + \underbrace{o(\|z_k - x^*\|)}_{=o(t_k)}\Big] \quad \text{①}$$

$$= \nabla c_i(x^*)^\top d + \underbrace{\frac{o(t_k)}{t_k}}_{\to 0}. \quad \text{so} \quad \nabla c_i(x^*)^\top d = 0.$$

If $i \in A(x^*) \cap \mathcal{I}$,
then

$$0 \leq \frac{1}{t_k} c_i(z_k) = \nabla c_i(x^*)^\top d + \underbrace{\frac{o(t_k)}{t_k}}_{\to 0}. \quad \text{so} \quad \nabla c_i(x^*)^\top d \geq 0.$$

$\underset{\substack{\uparrow \\ \text{Similar} \\ \text{to ①}}}{}$

Idea for the proof of (ii)

For $d \in \mathcal{F}(x^*)$ we need to find $z_k \underset{\underset{\vee}{0}}{\overset{\in \Omega}{,}} t_k \overset{\nearrow 0}{\underset{\vee}{,}}$ s.t.

$$\frac{z_k - x^*}{t_k} \to d \quad - (I)$$

choosing $z_k = x^* + t_k d$ works if $d$ points towards the interior of $\Omega$.

But doesn't work if $d$ is tangent to one of the level surface

If $c_i(z_k) = t_k \nabla c_i^T(x^*) d \quad \forall i \in \mathcal{A}(x^*)$ $\Big\}$ a system of $m = |\mathcal{A}(x^*)|$  $\underset{c_i(x)=0}{}$
then $z_k \in \Omega$. (why?) nonlinear equations in $n$ (*) variables.

On the other hand, if $z_k, t_k$ satisfy (I), then
$$c_i(z_k) \approx \underset{\underset{\approx 0}{}}{c_i(x^*)} + \nabla c_i(x^*)^T \underbrace{(z_k - x^*)}_{\approx t_k d} \approx t_k \nabla c_i(x^*)^T d$$

So the problem is essentially about solving the $\underset{\underset{\vee}{under\text{-}determined}}{}$ nonlinear system (*).

Generically, there should be $(n-m)$ d.o.f. in choosing the solution $z_k$ for a fixed $t_k$.



$$c_i(z_k) \approx \nabla c_i(x^*)^T (z_k - x^*)$$

not okay to choose $z_k = x^* + t_k d$,
   (as this may make $z_k$ infeasible.)
but it should be okay to choose $z_k \in \Omega$
so that
$$z_k - (x^* + t_k d) \ // \ \nabla c_i(x^*)$$
$$\Longleftrightarrow$$
$$P(z_k - (x^* - t_k d)) = 0 \quad\text{—— (**)}$$

$\nwarrow$ ortho-projection onto $\bigcap\limits_{i=1}^{m} \nabla c_i(x^*)^{\perp}$

$$\text{null}\left( \underbrace{\begin{bmatrix} \nabla c_1(x^*)^T \\ \vdots \\ \nabla c_m(x^*)^T \end{bmatrix}}_{\substack{\mathbb{R}^n \to \mathbb{R}^m \\ \text{full rank}}} \right)$$
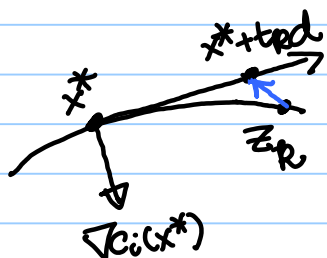
So, pick any basis $b_1, -, b_{n-m}$ of the null space, and set
$$Z^T = \begin{bmatrix} b_1^T \\ \vdots \\ b_{n-m}^T \end{bmatrix}.$$
$(n-m) \times n$

Then (**) $\Longleftrightarrow$ $Z^T(z_k - x^* - t_k d) = 0.$

So we look for $z_k$ by solving the (square) system of nonlinear equations

$$(***) \quad \begin{cases} c_i(z) - t_k \nabla c_i(x_*)^T d = 0 \quad \leftarrow \text{n vars, m eqts} \\ Z^T( z - x^* - t_k d) = 0 \quad \leftarrow \text{(n-m) extra eqts to remove redundancy} \end{cases}$$

for every small $t_k > 0$.

Of course, $(***)$ isn't truely necessary, an obique projection should also work.



It is fine to replace $Z^T$ above by any $B \in \mathbb{R}^{(m-n)\times n}$ so that

$$\begin{bmatrix} - \nabla c_1(x^*)^T - \\ \vdots \\ - \nabla c_m(x^*)^T - \\ B \end{bmatrix} \text{ is invertible}$$

[End of the intuitive explanation of the proof, now the rigorous proof:]

Proof of (ii). For notational convenience, assume $c_1, \cdots, c_m$ are the active constraints, ie. $\mathcal{A}(x^*) = \{1, \cdots, m\}$.

Write $C(z) = \begin{bmatrix} c_1(z) \\ \vdots \\ c_m(z) \end{bmatrix}$, $A(z^*) = \begin{bmatrix} \nabla c_1(x^*)^T \\ \vdots \\ \nabla c_m(x^*)^T \end{bmatrix} \in \mathbb{R}^{m \times n}$. ($m \leq n$ as the gradients are linearly indep.)

Now, assume $d \in \mathcal{F}(x^*)$, ie. $\nabla c_i^T(x^*)d = 0$, $i \in \mathcal{E}$
$\nabla c_i^T(x^*)d \geq 0$, $i \in \mathcal{I}$

Let $t_k > 0$ be s.t. $t_k \to 0$. Our goal is to find $z_k \in \Omega$ s.t. $\frac{z_k - x^*}{t_k} \to d$.

(C)

Let $B \in \mathbb{R}^{(n-m) \times m}$ so that $\begin{bmatrix} A(x^*) \\ B \end{bmatrix} \in \mathbb{R}^{n \times n}$ is non-singular, and consider the parametrized system of equations

$(\bigstar)$  $R(z, t) := \begin{bmatrix} C(z) - t A(x^*)d \\ B(z - x^* - td) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$, $R: \mathbb{R}^{n+1} \to \mathbb{R}^n$.

Claim: the solutions $z = z_k$ of this system for small $t = t_k > 0$ give a feasible sequence that approaches $x^*$ and satisfies the condition (C).

Note: $\nabla_z R(x^*, 0) = \nabla_z \begin{bmatrix} C(z) \\ B(z-x^*) \end{bmatrix} = \begin{bmatrix} A(x^*) \\ B \end{bmatrix}$ ← invertible

So, by the implicit function theorem, $\exists \varepsilon > 0$ st $\forall t \in (-\varepsilon, \varepsilon)$, there is a unique solution $z(t)$ of the system $(\cancel{\star})$, ie.
$$R(z(t), t) = 0.$$

The solutions $z_k = z(t_k)$ are what we need! Check:

- $i \in \mathcal{E} \implies C_i(z_k) = t_k \underbrace{\nabla c_i(x^*)^T d}_{=0} = 0$    <span style="color:red">$\left[\begin{array}{l} i \notin \mathcal{A}(x^*), \; C_i(x^*) > 0 \\ z_k \to x^*, \text{ so } C_i(z_k) > 0 \\ \text{for large } k. \end{array}\right]$</span>

   $i \in \mathcal{I} \cap \mathcal{A}(x^*) \implies C_i(z_k) = \underbrace{t_k}_{>0} \underbrace{\nabla c_i(x^*)^T d}_{\geq 0} \geq 0$.    <span style="color:red">So $z_k$ is indeed feasible.</span>

- $0 = R(z_k, t_k) = \begin{bmatrix} C(z_k) - t_k A(x^*) d \\ B(z_k - x^* - t_k d) \end{bmatrix}$

    $= \begin{bmatrix} \overset{=0}{C(x^*)} + \overset{=A(x^*)}{\nabla c(x^*)^T}(z_k - x^*) + o(\|z_k - x^*\|) - t_k A(x^*) d \\ B(z_k - x^* - t_k d) \end{bmatrix}$

<span style="color:blue">invertible →</span> $= \begin{bmatrix} A(x^*) \\ B \end{bmatrix}(z - x^* - t_k d) + o(\|z_k - x^*\|)$

so
$$\frac{z - x^*}{t_k} = d + o\left(\frac{\|z_k - x^*\|}{t_k}\right)$$

← This relation says
that $\|q_k\|$ has to
be bounded.
But then
$o(\|q_k\|) \to 0.$

call it $q_k$

$\underbrace{\phantom{\frac{\|z_k - x^*\|}{t_k}}}$ $\|q_k\|$

we have $\dfrac{z - x^*}{t_k} \to d$ as $k \to \infty$.

Q.E.D.

Note: the LICQ condition can be dispensed with for linear constraints,
and can be replaced by a weaker condition in general.

more about these later. And I'll show you a fun example in the HW.
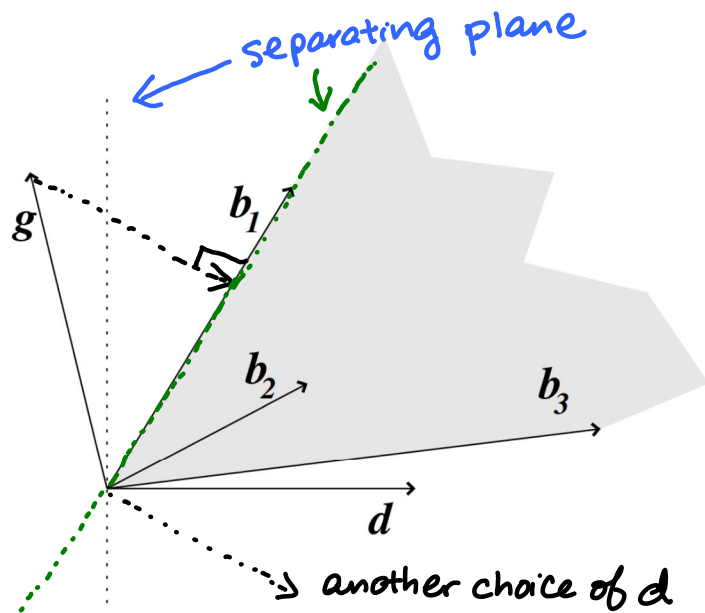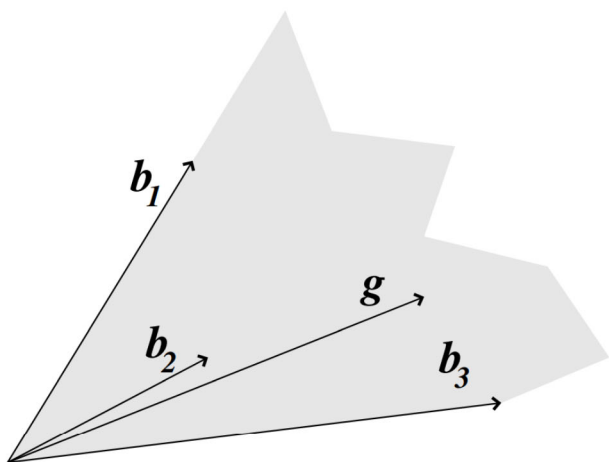
Step III

Lemma (Farkas)

Given any vectors $g$, $b_1, \cdots, b_m$, $c_1, \cdots, c_p \in \mathbb{R}^n$.
We have either:

- $g \in K = \{ By + Cw : y \geqslant 0 \}$, $B = [b_1, \cdots, b_m]$ $\quad^{n \times m}$ $\quad^{n \times p}$
  $C = [c_1, \cdots, c_p]$

OR else

- $\exists \, d \in \mathbb{R}^n$ st $g^T d < 0$, $B^T d \geqslant 0$, $C^T d = 0$.

But not both.



separating plane

For any such separating plane $d^\perp$,

$$g^T d < 0$$
$$s^T d \geqslant 0 \quad \forall s \in K$$
$$\Rightarrow (By + Cw)^T d \geqslant 0 \quad \forall y^{\geqslant 0}, w$$
$$= y^T B^T d + w^T C^T d$$
$$\Rightarrow B^T d \geqslant 0, \; C^T d = 0$$

So all we need is the existence of a separating plane, which seems obvious!

$b_1$ $\quad g$ $\quad b_2$ $\quad b_3$

$g$ $\quad b_1$ $\quad b_2$ $\quad b_3$ $\quad d$

→ another choice of $d$

Put differently, $\cancel{\exists} \, d \in \mathbb{R}^n$ st $g^T d < 0$, $B^T d \geqslant 0$, $C^T d = 0$ $\iff$ $g \in K$.

Proof (modulo a technical step):

(Easy step) We first show that the two alternatives cannot hold simultaneously.

If $g \in K$, i.e.
$$g = By + Cw \quad , \quad y \geqslant 0,$$
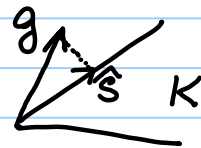and also
$$g^T d < 0, \quad B^T d \geqslant 0, \quad C^T d = 0$$
then
$$0 > g^T d = (By + Cw)^T d = \underbrace{y^T}_{\geqslant 0} \underbrace{B^T d}_{\geqslant 0} + w^T \underbrace{C^T d}_{= 0} \geqslant 0, \text{ a contradiction.}$$

(Harder step):

Assume $g \notin K$, we construct $d$ satisfying the properties.
We choose $d$ in the following way:

$$\text{Let } \hat{s} \in \underset{s \in K}{\text{argmin}} \, \|s - g\|_2^2, \quad \text{and} \quad d = \hat{s} - g.$$

<span style="color:red">A minimizer exists because <u>K is closed</u>,
but it requires some care to prove it. (We omit the argument here.
see Beck Ch6 or N-W Lemma 12.15)</span>

[It should be intuitively clear that $d^\perp$ separates $g$ from $K$. Here is a proof:]

It should be clear that $\hat{s}^T \perp \hat{s}-g$. (If not, it is because you forgot about the "conspiracy" between length and angle: $\|x\|_2^2 = \langle x, x \rangle = x^T x$.)
Precisely,

$\qquad$ $K$ is a cone, so $\alpha\hat{s} \in K$ $\forall \alpha \geq 0$. So $\|\alpha\hat{s}-g\|_2^2$ is minimized by $\alpha=1$.

So

$$\underbrace{\frac{d}{d\alpha}\|\alpha\hat{s}-g\|_2^2}_{\alpha^2\hat{s}^T\hat{s}-2\alpha\hat{s}^Tg+g^Tg}\Big|_{\alpha=1} = 0 \iff 2\alpha\hat{s}^T\hat{s} - 2\hat{s}^Tg\Big|_{\alpha=1} = 0 \iff \hat{s}^T(\hat{s}-g) = 0.$$
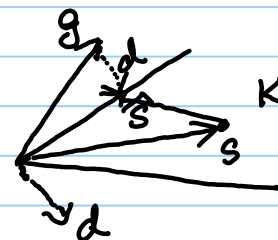
Now, note that K is not just a cone, it is also convex. So:



$$\|\underbrace{(1-\theta)\hat{s} + \theta s}_{=\hat{s}+\theta(s-\hat{s})} - g\|_2^2 \geq \|\hat{s}-g\|_2^2 \qquad \forall s \in K, \theta \in [0,1].$$

$$\Rightarrow \quad \langle \hat{s}-g + \theta(s-\hat{s}), \hat{s}-g+\theta(s-\hat{s})\rangle \geq \langle \hat{s}-g, \hat{s}-g\rangle$$

$$\Rightarrow \quad 2\theta(s-\hat{s})^T(\hat{s}-g) + \theta^2\|s-\hat{s}\|_2^2 \geq 0 \xrightarrow{\theta \downarrow 0} (s-\hat{s})^T(\underbrace{\hat{s}-g}_{d}) \geq 0$$

$$\text{so} \quad s^T d - \underbrace{\hat{s}^T d}_{=0} \geq 0, \text{ or } s^T d \geq 0 \ \forall s \in K$$

Also $d \neq 0$ since $g \notin K$, so $d^T g = d^T(\hat{s}-d) = 0 - d^T d = -\|d\|_2^2 < 0$.

We have shown that $d^1$ separates $K$ from $g$.

So $\quad d^T(By + Cw) \geqslant 0 \quad \forall \, y^{\geqslant 0}, w$

$\quad$ set $\, w=0, \, \underbrace{d^T B y}_{(B^T d)^T y} \geqslant 0 \, \forall y \geqslant 0$. This is only possible if $B^T d \geqslant 0$

Similarly, set $y=0, \, \underbrace{d^T C w}_{(C^T d)^T} = 0 \, \forall w$. This is only possible if $C^T d = 0$.

$\hspace{12cm}$ Q.E.D.

Comments on the proof:

(i) That $K = \{Cy + Bw : y \geqslant 0\}$ is closed is essential for the proof.

(ii) That $K$ is convex is very essential for the existence of a separating hyperplane. <span style="color:red">(The existence has little to do with the fact that $K$ is a closed cone.)</span> This is also why I like the treatment of Farkas' lemma in [Becks] more, except that it is longer.

$g \bullet$

$\uparrow$ convex

$\leftarrow$ not convex

(iii) We can state Farkas' lemma in a slightly simpler form without losing generality:

---

**Lemma (Farkas)**

Given any vectors $g$, $b_1, \cdots, b_m$, ~~$c_1, \cdots, c_p$~~ $\in \mathbb{R}^n$.
We have either:

- $g \in K = \{ \overset{n \times m}{B}y \pm \overset{n \times p}{Cw} : y \geqslant 0 \}$, $B = [b_1, \cdots, b_m]$

$\qquad\qquad$ OR else $\qquad\qquad$ ~~$C = [c_1, \cdots, c_p]$~~

- $\exists \, d \in \mathbb{R}^n$ st $g^T d < 0$, $B^T d \geqslant 0$, ~~$C^T d = 0$~~.

But not both.

---

Why doesn't it lose generality?

It is because we can always write a real number as the difference of two non-negative numbers, so

$$K = \{ By + Cw : y \geqslant 0 \} = \left\{ [B, C, -C] \begin{bmatrix} y \\ w^+ \\ w^- \end{bmatrix} : \begin{bmatrix} y \\ w^+ \\ w^- \end{bmatrix} \geqslant 0 \right\}.$$

$\underbrace{\qquad\qquad}$

↑ call this the new B

↑ the new y

# Proof of the KKT theorem:

If $x^*$ is a local solution, the LICQ is satisfied at $x^*$, then

$$\nabla f(x^*)^T d \geq 0, \quad \forall d \in T_{lin}(x^*) = F(x^*)$$

So by Farkas' lemma,

$$\nabla f(x^*) = \sum_{i \in A(x^*)} \lambda_i \nabla c_i(x^*), \quad \text{for some Lagrange multipliers } \lambda_i$$

$$\lambda_i \geq 0 \text{ for } i \in I \cap A(x^*)$$

To complete the proof, all we need is to define the vector $\lambda^*$ by

$$\lambda^*_i = \begin{cases} \lambda_i & i \in A(x^*) \\ 0 & i \in I \setminus A(x^*). \end{cases}$$

The rest is trivial to check.                    Q.E.D.