



262B-Lecture 16

Date created: 2021.03.16
N. of Pages: 17

$$\min f(x) + h(x)$$

$$\text{s.t. } x \in X$$

, assume:

f : convex & differentiable

h : convex but maybe non-smooth

Proximal gradient algorithm:

$$\begin{cases} z^{(k)} = x^{(k)} - \alpha^{(k)} \nabla f(x^{(k)}) & \rightarrow \text{gradient alg. on } f(\cdot) \\ x^{(k+1)} = \text{prox}_{\alpha^{(k)}, h}(z^{(k)}) & \rightarrow \text{proximal on } h(\cdot) \text{ and } X \end{cases}$$

$$\text{argmin}_{x \in X} \left(h(x) + \frac{1}{2\alpha^{(k)}} \|x - z^{(k)}\|^2 \right)$$

Two special cases:

$$1 - f(x) = 0$$

$$\Rightarrow z^{(k)} = x^{(k)}$$

$$\Rightarrow x^{(k+1)} = \text{prox}_{\alpha^{(k)}, h}(x^{(k)})$$

\Rightarrow proximal algorithm for $\min h(x)$ s.t. $x \in X$

$$2 - h(x) = 0$$

$$\text{prox}_{\alpha^{(k)}, h}(z) = \text{argmin}_{x \in X} \left(0 + \frac{1}{2\alpha^{(k)}} \|x - z\|^2 \right)$$

$$= \text{argmin}_{x \in X} \|x - z\| = \mathcal{P}_X(z)$$

$$\Rightarrow x^{(k+1)} = \text{prox}_{\alpha^{(k)}, h} (z^{(k)}) = \mathcal{P}_X (z^{(k)})$$

$$= \mathcal{P}_X (x^{(k)} - \alpha^{(k)} \nabla f(x^{(k)}))$$

\Rightarrow Proximal gradient method is a generalization of projected gradient method.

Ex: $\min_{x \in \mathbb{R}^n} f(x) + |x|$, \rightarrow sparse solution in machine learning

special case: Lasso

$$\min \|Ax - b\|^2 + \lambda |x| \rightarrow \min \frac{1}{\lambda} \|Ax - b\|^2 + |x|$$

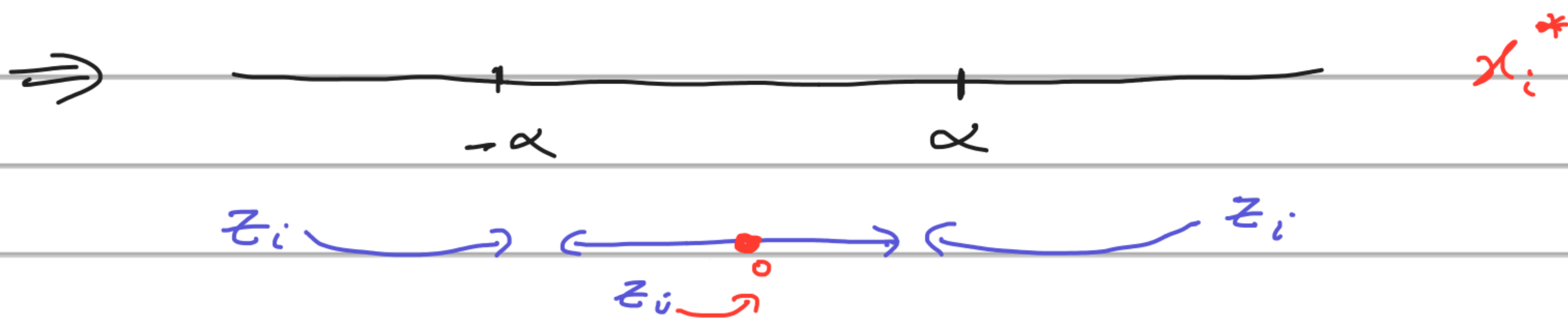
proximal operator: $h(x) = |x|$

$$\text{prox}_{\alpha, h} (z) \rightarrow \min_x |x| + \frac{1}{2\alpha} \|x - z\|^2$$

$$\rightarrow \min_x \sum_{i=1}^n (|x_i| + \frac{1}{2\alpha} (x_i - z_i)^2)$$

decomposable

$$\Rightarrow x_i^* = \begin{cases} z_i - \alpha & \text{if } z_i > \alpha \\ 0 & \text{if } z_i \in [-\alpha, \alpha] \\ z_i + \alpha & \text{if } z_i < -\alpha \end{cases}$$



Prox operator in this case is called soft thresholding or shrinkage operator.

$$\min_{x \in \mathbb{R}^n} f(x) + |x|,$$

$$\Rightarrow x^{(k+1)} = \underbrace{ST}_{\alpha^{(k)}} (x^{(k)} - \alpha^{(k)} \nabla f(x^{(k)}))$$

soft thresholding

Recall : - proximal gradient algorithm is a general case of projected gradient algorithm.

- Convergence proofs for projected gradient algorithms are based on :

$$\|P_X(x) - P_X(y)\| \leq \|x - y\| \Rightarrow \text{non-expansive property}$$

$$\text{Thm: } \|\text{prox}_{\alpha, h}(x) - \text{prox}_{\alpha, h}(y)\| \leq \|x - y\| \quad \forall x, y \in \mathbb{R}^n$$

proof in a special case: $h(\cdot)$ = differentiable

shorthand notation: $z(x) = \text{prox}_{\alpha, h}(x)$

$$z(y) = \text{prox}_{\alpha, h}(y)$$

$$\Rightarrow z(x) = \underset{w \in X}{\text{argmin}} \left(h(w) + \underbrace{\frac{1}{2\alpha} \|w - x\|^2}_{\text{convex}} \right)$$

Foc:

gradient of objective =

$$\nabla h(z(x)) + \frac{z(x) - x}{\alpha}$$

$$\left(\nabla h(z(x)) + \frac{z(x) - x}{\alpha} \right)^T (w - z(x)) \geq 0$$

$\forall w \in X$

since $z(y) \in X$, pick $w = z(y)$

$$\Rightarrow \left(\nabla h(z(x)) + \frac{z(x) - x}{\alpha} \right)^T (z(y) - z(x)) \geq 0$$

①

Redo for $z(y)$:

$$\left(\nabla h(z(y)) + \frac{z(y) - y}{\alpha} \right)^T (z(x) - z(y)) \geq 0$$

(2)

Also, convexity of $h(\cdot)$:

$$\left(\nabla h(z(x)) - \nabla h(z(y)) \right)^T (z(x) - z(y)) \geq 0$$

(3)

(1), (2), (3) :

$$\left((z(x) - x) - (z(y) - y) \right)^T (z(x) - z(y)) \leq 0$$

$$\Rightarrow \|z(x) - z(y)\|^2 \leq (x - y)^T (z(x) - z(y))$$

$$\leq \|x - y\| \times \|z(x) - z(y)\|$$

(*)

1 - If $\|z(x) - z(y)\| = 0 \Rightarrow \|z(x) - z(y)\| \leq \|x - y\|$

2 - If $\|z(x) - z(y)\| \neq 0$

$$\Rightarrow (*) : \|z(x) - z(y)\| \leq \|x - y\|$$

Similar to projected gradient method:

$$1 - \|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\| \quad (\forall x, y \in \mathbb{R}^n)$$

$$\Rightarrow \text{iteration complexity} = O\left(\frac{1}{\varepsilon}\right)$$

$$2 - \|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\| \quad \forall x, y \in \mathbb{R}^n$$

$$\nabla^2 f(x) \succeq m I, \quad m > 0, \quad \forall x \in \mathbb{R}^n$$

$$\Rightarrow \|x^{(k+1)} - x_*\| \leq \sqrt{1 - \frac{m}{L}} \|x^{(k)} - x_*\|$$

$$\text{if } \alpha^{(k)} = \frac{1}{L}$$

$$\Rightarrow \text{iteration complexity} = O\left(\log \frac{1}{\varepsilon}\right)$$

Note: Number of iterations is like

unconstrained case, but we should take

a proximal step at every iteration:

per-iteration complexity $\sim h(\cdot), X$

Chapter 4:

previously, we studied:

$$\left\{ \begin{array}{l} \min f(x) \longrightarrow \text{non-convex} \\ \text{s.t. } x \in X \longrightarrow \text{convex} \longrightarrow \text{geometric FOC, SOC} \end{array} \right.$$

What if $X = \text{non-convex}$?

- optimization with equality constraints:

$$\min f(x)$$

$$x \in \mathbb{R}^n$$

$$\text{s.t. } \left. \begin{array}{l} h_1(x) = 0 \\ h_2(x) = 0 \\ \vdots \\ h_m(x) = 0 \end{array} \right\} \rightarrow \begin{array}{l} h(x) = 0 \\ \text{vector} \end{array} \rightarrow \text{nonlinear}$$

Definition: Cone of first-order feasible

directions at

$$V(y) = \left\{ \Delta x \mid \nabla h_i(y)^T \Delta x = 0, \quad i = 1, \dots, m \right\}, \quad y \in \mathbb{R}^n$$

variation

Idea: pick a feasible y : $h(y) = 0$

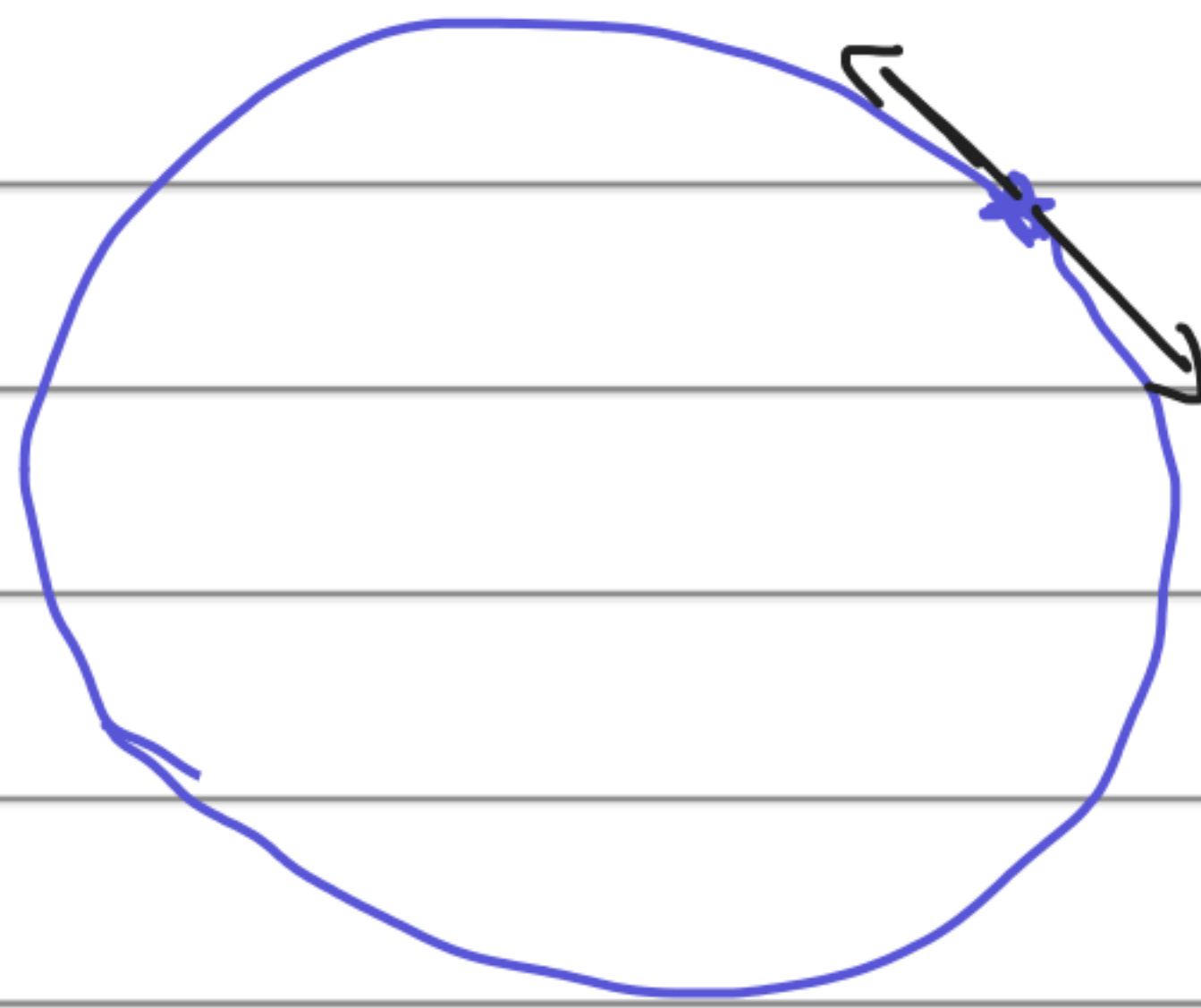
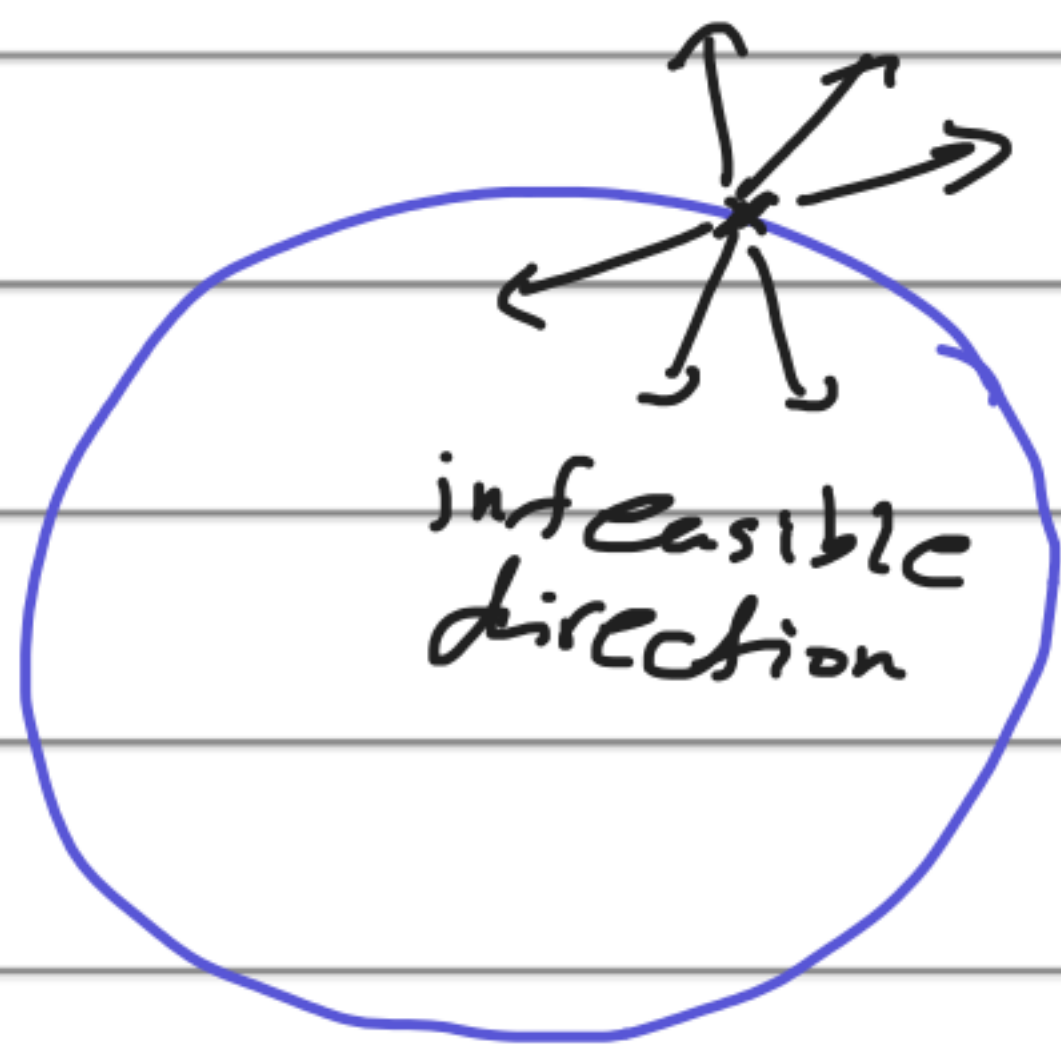
perturb it as $y + \Delta y$ s.t.: $h(y + \Delta y) = 0$

$$h_i(y) = 0, \quad 0 = h_i(y + \Delta y) = \cancel{h_i(y)} + \nabla h_i(y)^T \Delta y + \underbrace{O(\|\Delta y\|^2)}_{\text{small}}$$

If Δy : small $\Rightarrow \nabla h_i(y)^T \Delta y \approx 0$

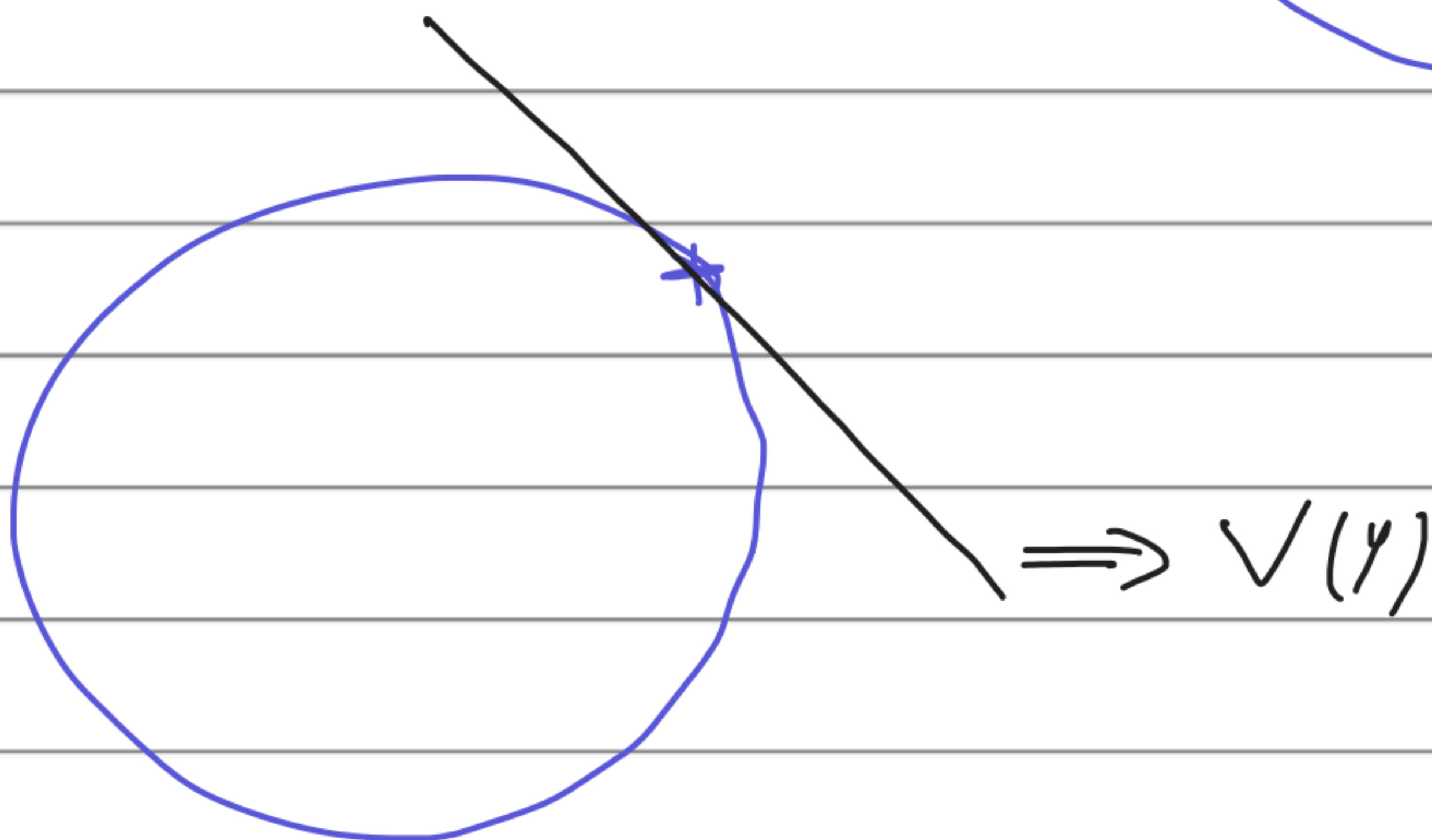
$\Rightarrow \Delta x \in V(y)$: $h(y + \underbrace{\varepsilon \Delta x}_{\Downarrow \text{almost feasible}}) \approx 0$ if $\varepsilon = \text{small}$

Ex:



not feasible,
but almost
feasible

\Rightarrow



$V(y)$ is the same as Tangent plane
at y if y is a regular point.

Def: y is a regular point if $\nabla h_1(y), \dots,$
 $\nabla h_m(y)$ are linearly independent.

FOC (necessary): If x_* is a regular point
and a local min for $\min f(x) \text{ s.t. } h(x) = 0,$
then $\exists \lambda_1^*, \dots, \lambda_m^* \in \mathbb{R}$ s.t.

$$\nabla f(x_*) + \sum_{i=1}^m \lambda_i^* \nabla h_i(x_*) = 0$$

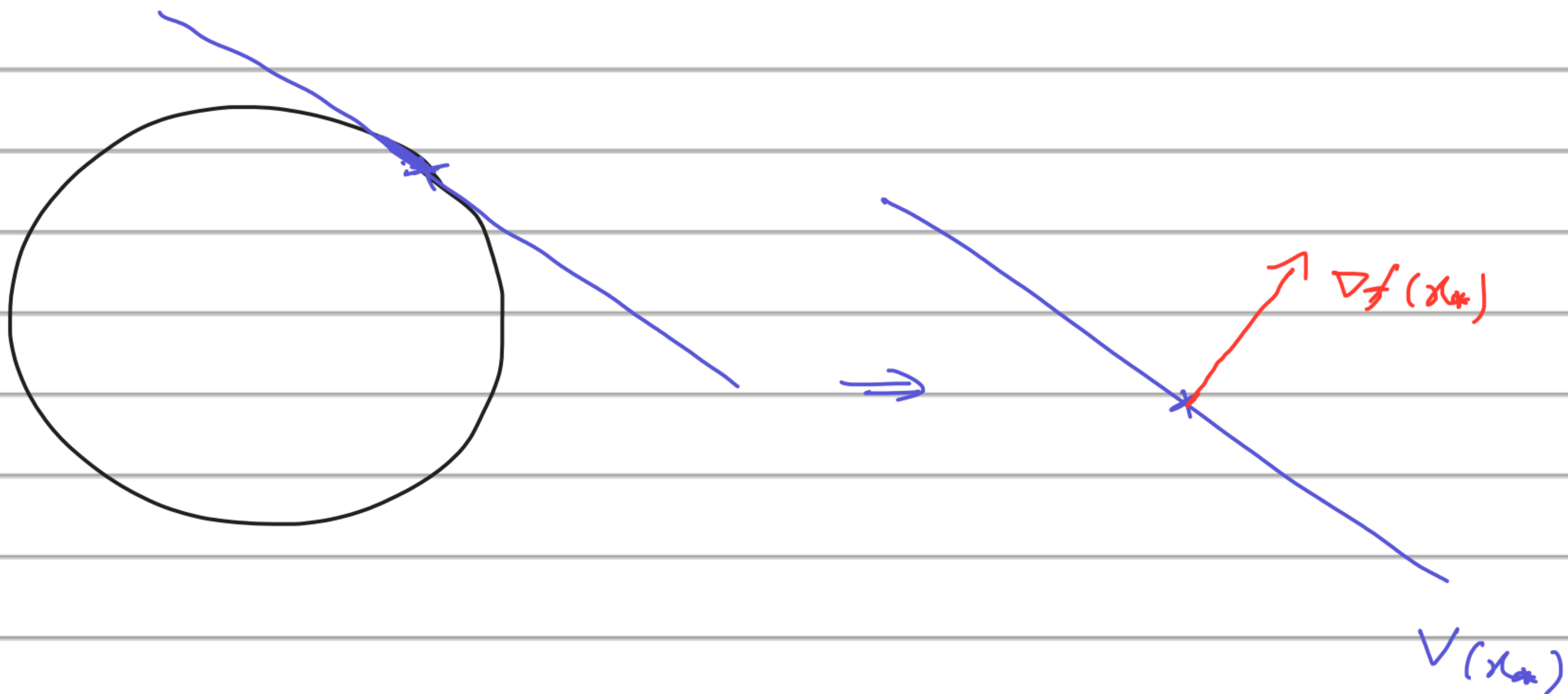
*

λ_i^* : Lagrange multipliers

Interpretations:

- 1 - $\nabla f(x_*) \in$ space spanned by gradients of constraints
- 2 - $\nabla f(x_*)$ is orthogonal to $V(x_*)$.

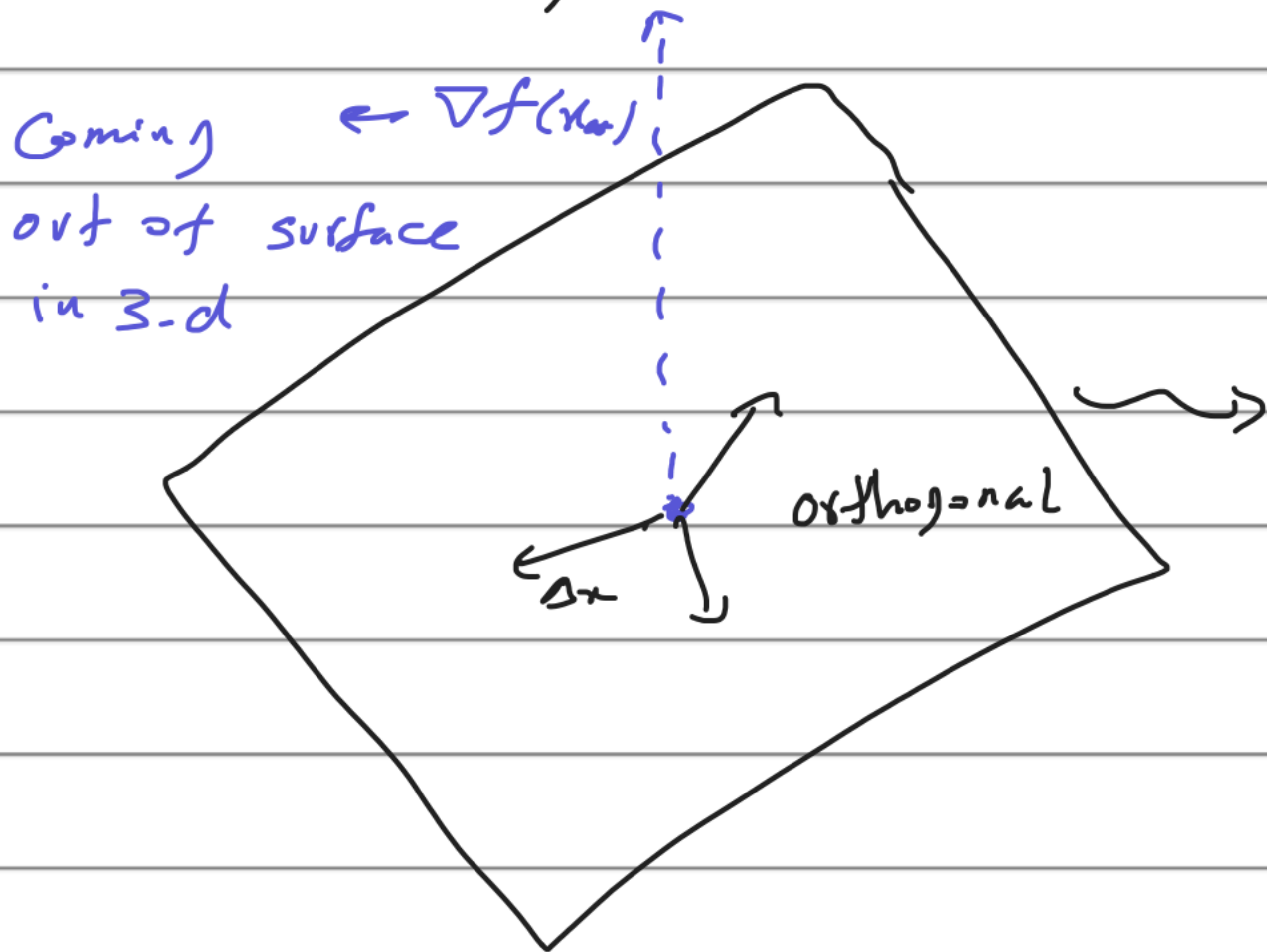
Ex:



Proof: $(*) \Rightarrow$ If Δx satisfies

$$\nabla h_i(x_*)^T \Delta x = 0 \quad i=1, \dots, m, \text{ then } \nabla f(x_*)^T \Delta x = 0$$

\Rightarrow Every $\Delta x \in V(x_*)$ is orthogonal to $\nabla f(x_*)$



surface defined by
gradients of constraints

Soc (necessary) : If x_* is a regular point and a local min, then $\lambda_1^*, \dots, \lambda_m^*$ satisfying

Foc also satisfy:

$$\Delta x^T \left(\nabla^2 f(x_*) + \sum_{i=1}^m \lambda_i^* \nabla^2 h_i(x_*) \right) \Delta x \geq 0$$
$$\forall \Delta x \in V(x_*)$$

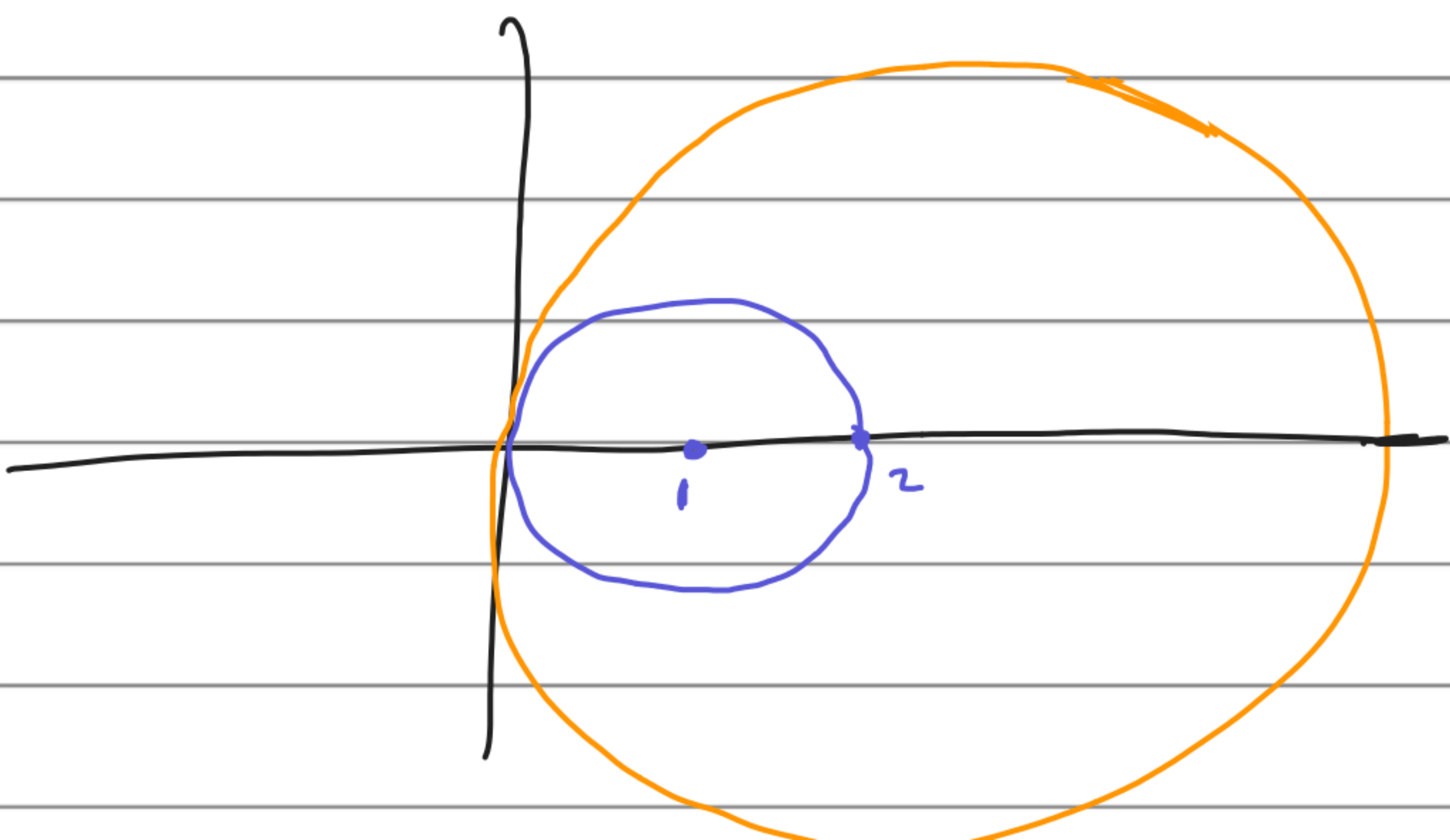
Why is regularity important:

$$\min x_1 + x_2$$

$$\text{s.t. } (x_1 - 1)^2 + x_2^2 - 1 = 0$$

$$(x_1 - 2)^2 + x_2^2 - 4 = 0$$

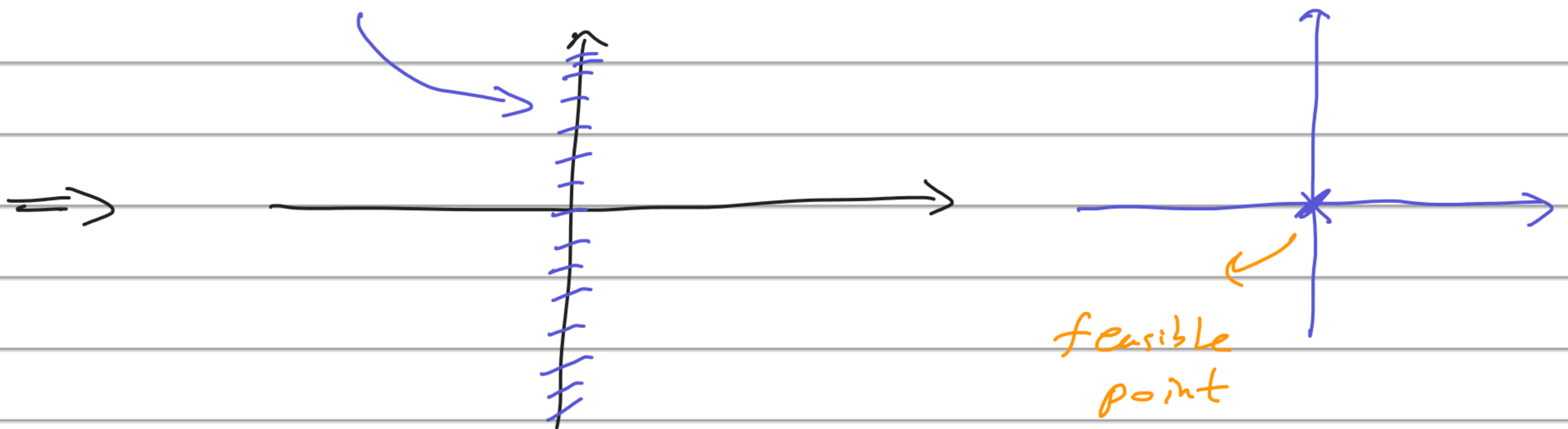
} \Rightarrow feasible set



x
origin = feasible
point

$$\Rightarrow x_* = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\text{Also, } \underbrace{V(x_*)}_{\text{tangent plane}} = \left\{ \Delta x \mid \begin{bmatrix} -2 & 0 \end{bmatrix} \Delta x = 0, \begin{bmatrix} -4 & 0 \end{bmatrix} \Delta x = 0 \right\}$$



Line is not tangent plane
for a single point

x_* : Not a regular point ($\begin{bmatrix} -2 \\ 0 \end{bmatrix}, \begin{bmatrix} -4 \\ 0 \end{bmatrix}$:
linearly dependent)

$$\nabla f(x_*) = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\text{FOC: } \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \lambda_1^* \begin{bmatrix} -2 \\ 0 \end{bmatrix} + \lambda_2^* \begin{bmatrix} -4 \\ 0 \end{bmatrix} = 0$$

This has no solution.

Prove FOC, SOC:

- 1 - penalty method
- 2 - Elimination method

1 - penalty method:

Define:
$$F^K(x) = \underbrace{f(x)}_{\text{surrogate for original optimization}} + \frac{K}{2} \underbrace{\|h(x)\|^2}_{\text{constraint violation}} + \frac{\alpha}{2} \underbrace{\|x - x_*\|^2}_{\text{proximal term}}$$

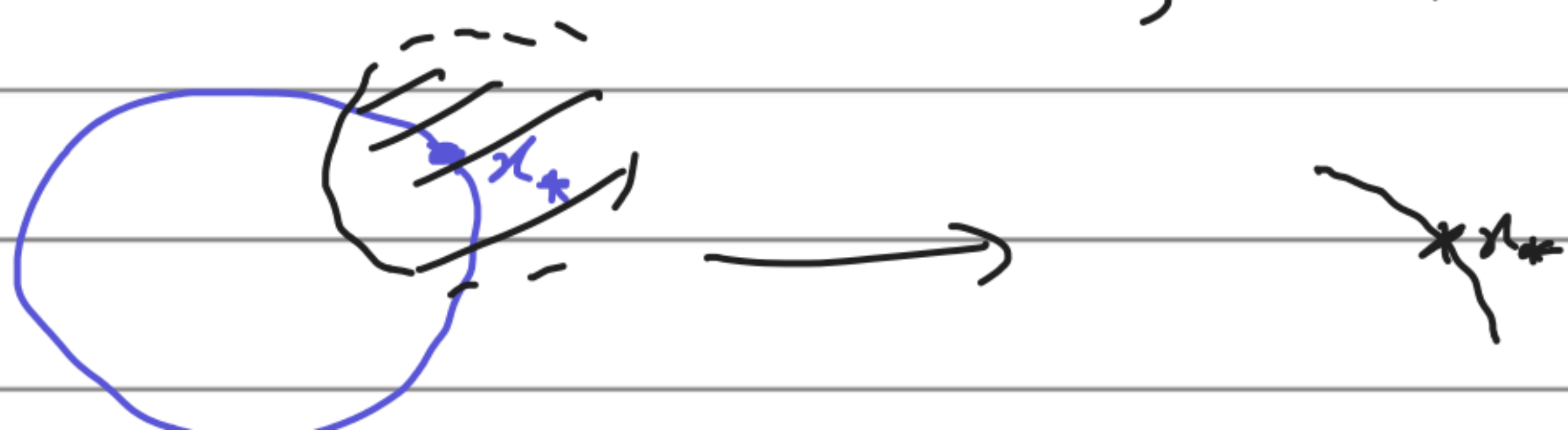
K : positive integer, $\alpha > 0$

$$F^1(x) \rightarrow F^2(x) \rightarrow \dots \rightarrow F^K(x) \rightarrow \dots$$

Since x_* is a local min, $\exists \epsilon > 0$ s.t.

$$f(x_*) \leq f(x) \quad \forall x \in \Sigma \text{ s.t. } h(x) = 0$$

where $\Sigma = \{x \mid \|x - x_*\| \leq \epsilon\}$



Define: $x^{(k)}$ = global min for

$$\begin{array}{l} \min F^k(x) \\ \text{s.t. } x \in S \end{array}$$

no $h(x)=0$ in this optimization

$$\text{Note: } F^k(x^{(k)}) \leq F^k(\underbrace{x_*}_{\text{arbitrary point in } S}) = f(x_*)$$

$$+ \frac{k}{2} \|h(x_*)\|^2 + \frac{\alpha}{2} \|x_* - x_*\|^2$$

$$\Rightarrow F^k(x^{(k)}) \leq f(x_*)$$

$$\Rightarrow \left(F^k(x^{(k)}) = f(x^{(k)}) + \frac{k}{2} \|h(x^{(k)})\|^2 + \frac{\alpha}{2} \|x^{(k)} - x_*\|^2 \right) \leq f(x_*) \quad (*)$$

$f(x^{(k)})$ is bounded on S , so the left

side would blow up and violate $(*)$

when $k \rightarrow \infty$, unless $\|h(x^{(k)})\| \rightarrow 0$

$$\Rightarrow \lim_{k \rightarrow \infty} \|h(x^{(k)})\| = 0$$

let \bar{x} denote an arbitrary limit point of

$$\{x^{(k)}\} \Rightarrow h(\bar{x}) = 0 \Rightarrow \bar{x} : \text{feasible point}$$

$$\textcircled{*} \Rightarrow f(x^{(k)}) + \frac{\alpha}{2} \|x^{(k)} - x_*\|^2 \leq f(x_*)$$

$$k \rightarrow \infty : f(\bar{x}) + \frac{\alpha}{2} \|\bar{x} - x_*\|^2 \leq f(x_*)$$

$$\text{If } \bar{x} \neq x_* \Rightarrow \|\bar{x} - x_*\| > 0 \Rightarrow$$

$$f(\bar{x}) < f(x_*) \rightarrow \text{Contradiction}$$

(x_* : best feasible point in the ball Σ)

$$\Rightarrow \bar{x} = x_* \rightarrow \text{every limit point is } x_*$$

$$\Rightarrow \begin{cases} \min F^k(x) \\ \text{s.t. } x \in \Sigma \end{cases}$$

or

$$\min f(x) + \frac{k}{2} \|h(x)\|^2 + \frac{\alpha}{2} \|x - x_*\|^2$$

$$\text{s.t. } x \in \Sigma$$

As $k \rightarrow \infty$, solution of new problem

converges x_* .

Since $x_* \in \text{interior of } S$

$\Rightarrow x^{(k)} \in \text{interior of } S \text{ if } k \text{ is large}$

$\Rightarrow \min F^k(x)$ if k is large.

~~s.t. $x \in S$~~

\Rightarrow Unconstrained optimization.

FOC, SOC for unconstrained optimization:

$$\nabla F^k(x^{(k)}) = 0, \quad \nabla^2 F^k(x^{(k)}) \succeq 0$$

Goal: \downarrow
= FOC for
constrained opt
as $k \rightarrow \infty$

\downarrow
= SOC for
constrained opt
as $k \rightarrow \infty$

$$\nabla h(x) = [\nabla h_1(x) \quad \nabla h_2(x) \quad \dots \quad \nabla h_m(x)]$$

regularity: $\nabla h(x_*) = \text{full column rank}$

$$\{x^{(k)}\} \rightarrow x_*$$

$$\Rightarrow \nabla h(x^{(k)}) = \text{full column rank if } k = \text{large}$$

$$\Rightarrow \nabla h(x^{(k)})^T \nabla h(x^{(k)}) = \text{invertible matrix}$$

$$\underbrace{0}_{\text{Foc}} = \nabla F^k(x^{(k)}) = \nabla f(x^{(k)}) + \underbrace{k}_{=} \nabla h(x^{(k)}) \underbrace{h(x^{(k)})}_{=} + \alpha(x^{(k)} - x_*)$$

$$\Rightarrow k h(x^{(k)}) = - (\nabla h(x^{(k)})^T \nabla h(x^{(k)}))^{-1} \nabla h(x^{(k)})^T (\nabla f(x^{(k)}) + \alpha(x^{(k)} - x_*))$$

$$k \rightarrow \infty : \lim_{k \rightarrow \infty} k h(x^{(k)}) = \lambda_*$$

$$\text{where: } \lambda_* = - (\nabla h(x_*)^T \nabla h(x_*))^{-1} \nabla h(x_*)^T \nabla f(x_*)$$

~~*~~ ~~*~~ \rightarrow Take Limit:

$$\nabla f(x_*) + \underbrace{\sum_{i=1}^n \nabla h_i(x_*) \lambda_i^*}_{\lambda_*} = 0$$

explicitly found Lagrange multipliers