

FODA L16

Gradient Descent
#2

functions $f: \mathbb{R}^d \rightarrow \mathbb{R}$

$$\nabla f: \mathbb{R}^d \rightarrow \mathbb{R}^d \quad f(\alpha) \quad \alpha = (\alpha_1, \dots, \alpha_d)$$

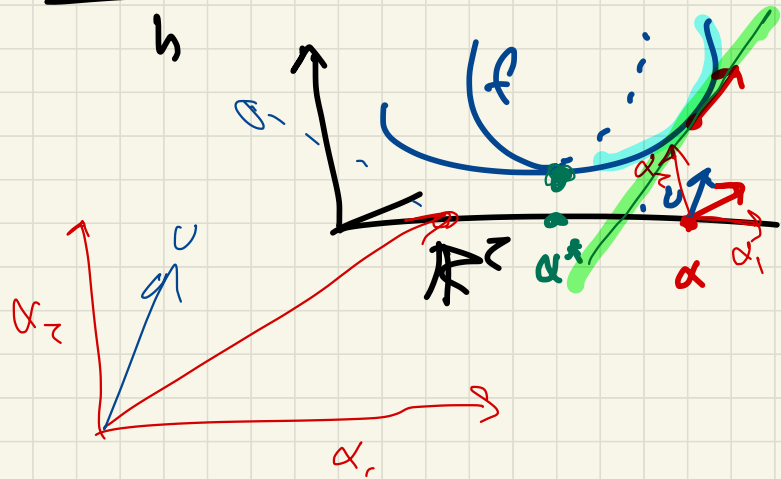
$$\nabla f = \left(\frac{df}{d\alpha_1}, \frac{df}{d\alpha_2}, \dots, \frac{df}{d\alpha_d} \right) \quad u = (u_1, u_2, \dots, u_d)$$

directional
derivative
 $\|u\|=1$

$$\nabla_u f(\alpha) = \lim_{h \rightarrow 0} \frac{f(\alpha + hu) - f(\alpha)}{h}$$

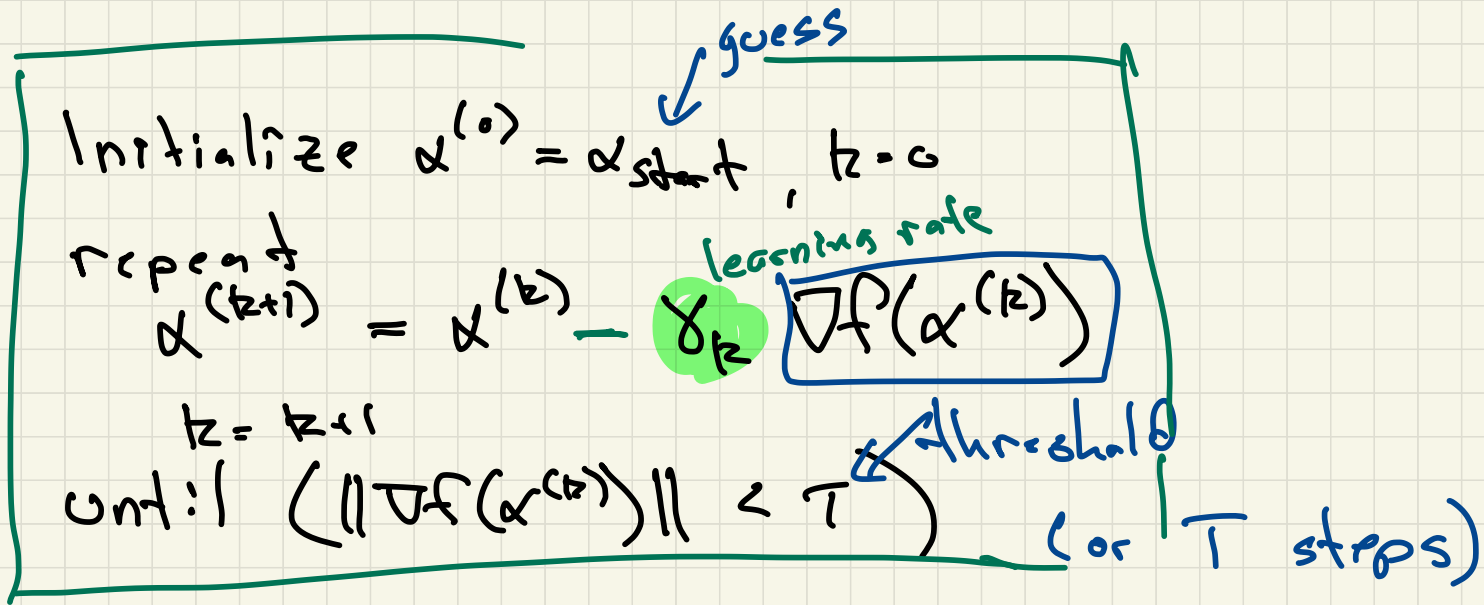
$$\nabla_u f(\alpha) = \langle \nabla f(\alpha), u \rangle$$

$$\max_{\substack{u \\ \|u\|=1}} \langle \nabla f(\alpha), u \rangle \rightarrow u = \frac{\nabla f(\alpha)}{\|\nabla f(\alpha)\|}$$



Gradient Descent Algo

Goal $\min_{\alpha \in \mathbb{R}^d} f(\alpha)$ or $\operatorname{argmin}_{\alpha \in \mathbb{R}^d} f(\alpha)$



Stopping condition

steps T

(finite constraint)

$$\|\nabla f(\alpha)\| \leq \epsilon$$

at optimal α^*

$$\nabla f(\alpha^*) = (0, 0, \dots, 0)$$

$$\|\nabla f(\alpha^*)\| = 0$$

Learning Rate γ how to choose?

Lipschitz function $g: \mathbb{R}^d \rightarrow \mathbb{R}^k$ is L -Lipschitz

if $\forall p, g \in \mathbb{R}^d$

$$\|g(p) - g(g)\| \leq L \overbrace{\|p - g\|}^h$$

$$u = \frac{p - g}{\|p - g\|}$$

$$g = p + uh$$

$$\frac{\|g(p) - g(p + uh)\|}{h} \leq L$$

Let $\nabla f = g$

Assume ∇f is L -Lipschitz

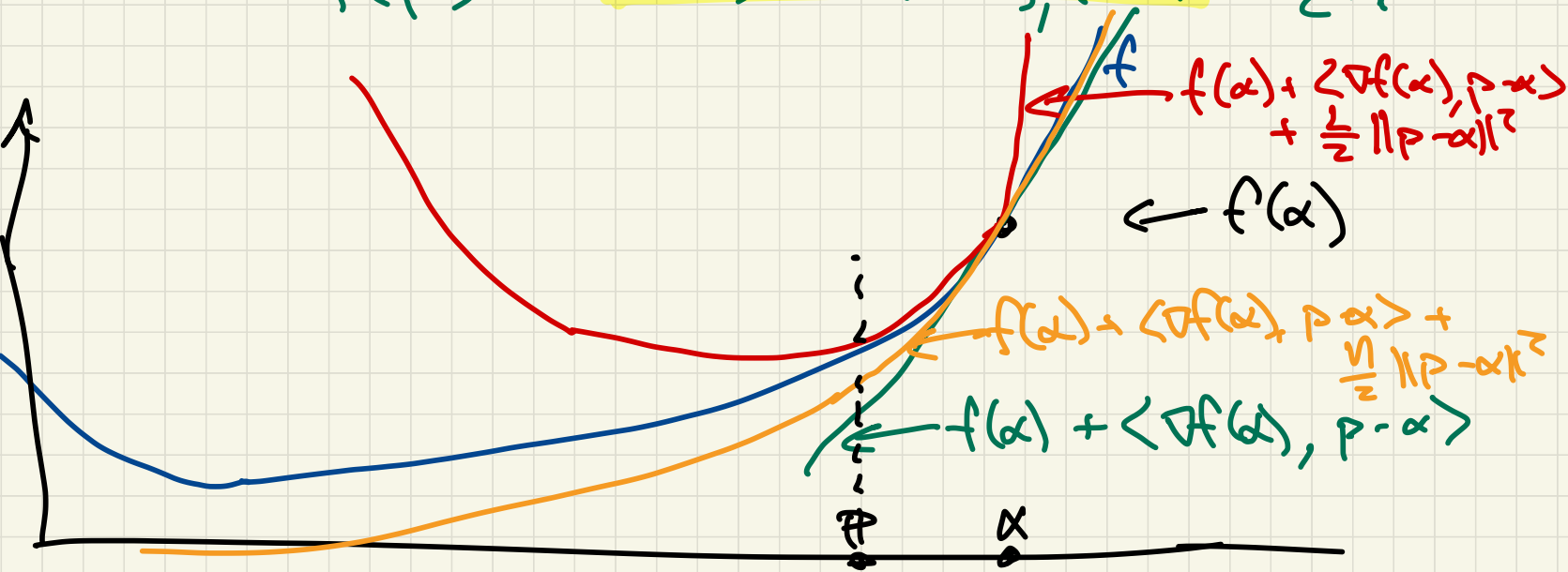
Set $\gamma < \frac{1}{L}$ then

GD will converge, if convex f
after $t = \frac{c/\epsilon}{\gamma}$ steps
 $f(x^{(t)}) - f(x^*) \leq \epsilon$

Strongly convex function (strict)

f η -strongly convex $\forall \alpha, p \in \mathbb{R}^d$

$$f(p) \geq f(\alpha) + \langle \nabla f(\alpha), p - \alpha \rangle + \frac{\eta}{2} \|p - \alpha\|^2$$



If f is η -strongly convex, \mathcal{D} of L -Lipschitz
set $\gamma \leq \frac{2}{(L+\eta)}$ after $k > C \cdot \log \frac{1}{\epsilon}$

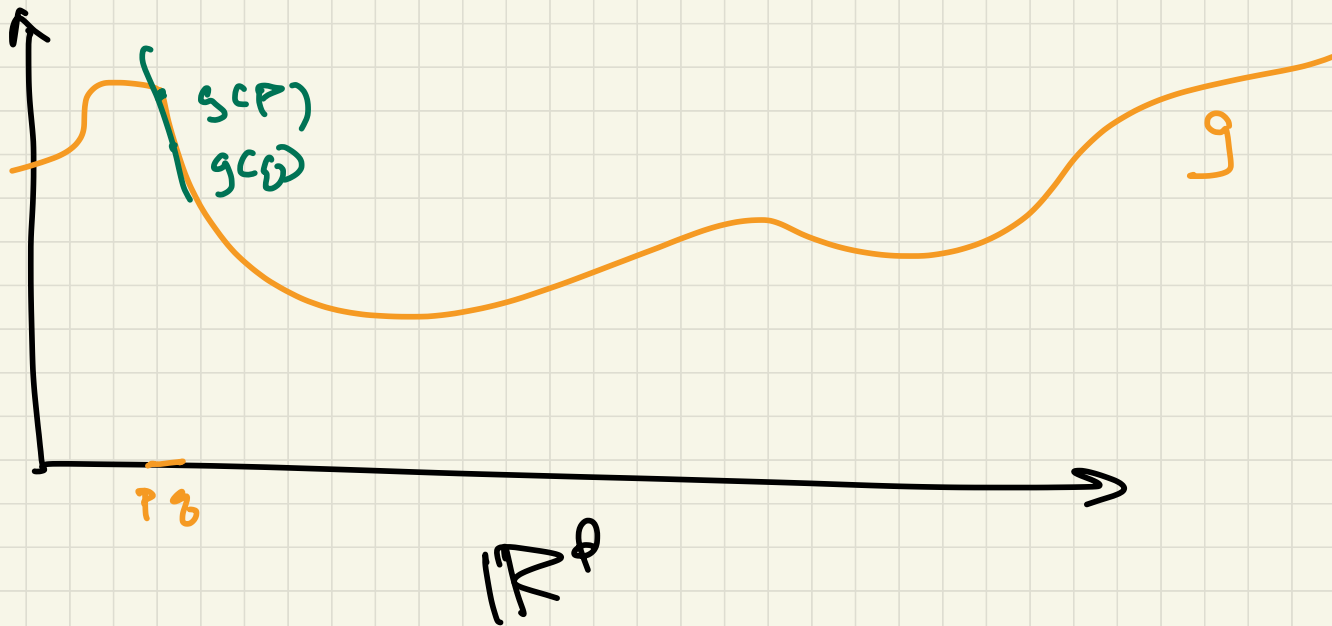
$$f(x^{(k)}) - f(x^*) \leq \epsilon$$

Imagine $C=1$, $\log = \log_{10}$

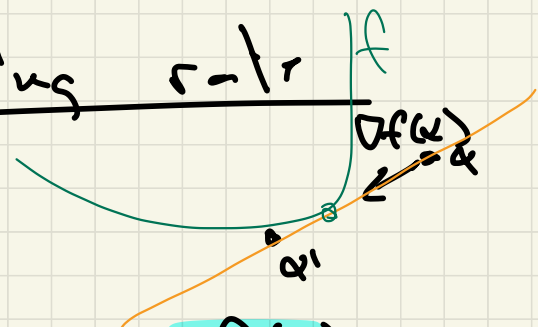
error $\epsilon = \underbrace{0.00 \dots 0}_{k} \times x \times r$

Lipschitz

g



Adjusting $\gamma \leftarrow$ learning rate



Line Search

Every step $\alpha' = \alpha - \gamma \nabla f(\alpha)$

goal minimize $f(\alpha')$

$$\gamma^* = \underset{\gamma \in \mathbb{R}}{\text{argmin}} f(\alpha - \gamma \nabla f(\alpha))$$

binary search

Often too slow

Adjustable Rate

Start w/ $\gamma = \gamma_0$

Sometimes $\gamma = \gamma \beta$

a bit larger than guess

$$\beta \in (0.1, 0.8)$$
$$\beta = 0.75$$

Condition

if

$$f(\alpha - \gamma \nabla f(\alpha)) > f(\alpha) - \frac{\gamma}{2} \|\nabla f(\alpha)\|^2$$

then $\gamma = \beta \gamma$

