

FODA - L17

Gradient Descent

SGD, on data

Gradient Descent $f: \mathbb{R}^d \rightarrow \mathbb{R}$

Goal: $\underset{\alpha \in \mathbb{R}^d}{\operatorname{arg\,min}} f(\alpha)$

Input $(X, y) = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\} \subset \mathbb{R}^d \times \mathbb{R}$

"loss" function $L((X, y), M_\alpha)$

$$f(\alpha) = L((X, y), M_\alpha) = \text{SSE}((X, y), M_\alpha)$$

$$= \sum_{(x_i, y_i) \in (X, y)} \left(M_{\alpha}(x_i) - y_i \right)^2$$

es. $\langle \alpha, x_i \rangle$

Model M_α for polynomial regression

$$\mathbb{R}^d = \mathbb{R}$$

$$x_i \in \mathbb{R}$$

$$p = 2$$

$$M_\alpha(x_i) = \langle \alpha, (1, x_i, x_i^2) \rangle = \alpha_0 + \alpha_1 x_i + \alpha_2 x_i^2$$
$$= \sum_{j=0}^2 \alpha_j x_i^j$$

$$\alpha = (\alpha_0, \alpha_1, \alpha_2) \in \mathbb{R}^3$$

$n=1$
single data point

$$\text{GD: } \alpha = \alpha - \gamma \nabla f(\alpha)$$

$$j = \{0, 1, 2\}$$

$$\frac{\partial}{\partial \alpha_j} f(\alpha) = \frac{\partial}{\partial \alpha_j} (M_\alpha(x_i) - y_i)^2 = 2(M_\alpha(x_i) - y_i) \frac{\partial}{\partial \alpha_j} (M_\alpha(x_i) - y_i)$$

$$= 2(M_\alpha(x_i) - y_i) \frac{\partial}{\partial \alpha_j} \left(\sum_{k=0}^2 \alpha_k x_i^k - y_i \right)$$

$$= 2(M_\alpha(x_i) - y_i) x_i^j$$

$$\frac{\partial}{\partial x_j} f(x) = 2 (M_x(x_j) - y_j) x_j^j \quad n=1$$

$$\begin{aligned} \nabla f(x) &= \left(\frac{\partial}{\partial x_0} f(x), \frac{\partial}{\partial x_1} f(x), \frac{\partial}{\partial x_2} f(x) \right) \\ &= 2 (M_x(x_j) - y_j) \cdot (1, x_j, x_j^2) \end{aligned}$$

LMS update rule, Widrow-Hoff update

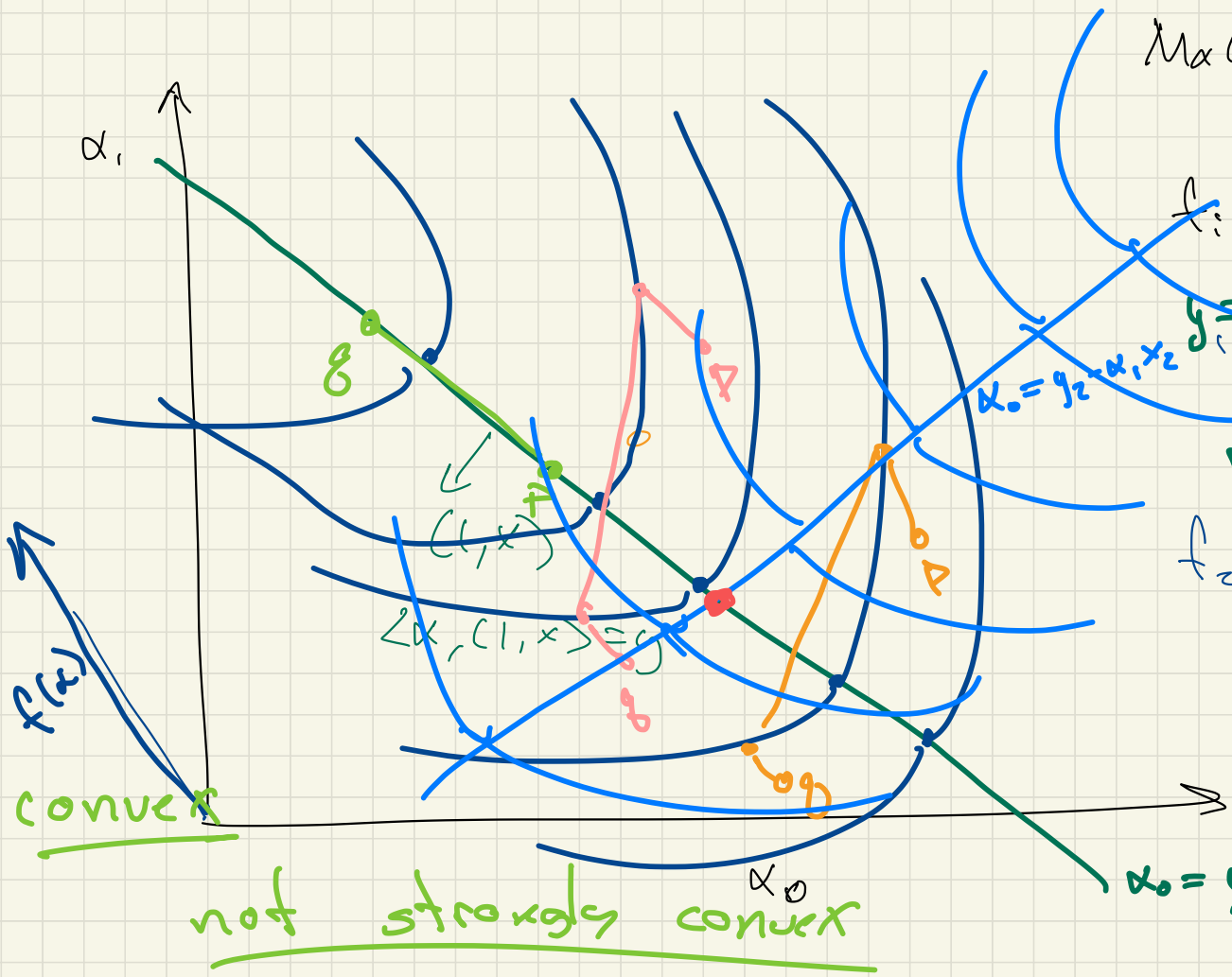
$$f(x) = \sum_{i=1}^n f_i(x)$$

$n > 1$

decomposable function

$$\begin{aligned} f_i(x) &= (Mx(x_i) - y_i)^2 \quad \leftarrow \text{each } f_i \\ &= (\alpha_0 + \alpha_1 x_i + \alpha_2 x_i^2 - y_i)^2 \quad \text{one observation } (x_i, y_i) \end{aligned}$$

if each f_i is ^(strongly) convex, then
 $f = \sum_i f_i$ is also ^(strongly) convex



$$M_{\alpha}(x_i) = \alpha_0 + \alpha_1 x_i$$

$$(x_i, y_i)$$

$$f_i(\alpha) = (y_i - \alpha_0 - \alpha_1 x_i)^2$$

$$y_i = \alpha_0 + \alpha_1 x_i$$

$$\Leftrightarrow f_i(\alpha) = 0$$

$$\alpha_0 = y_2 - \alpha_1 x_2$$

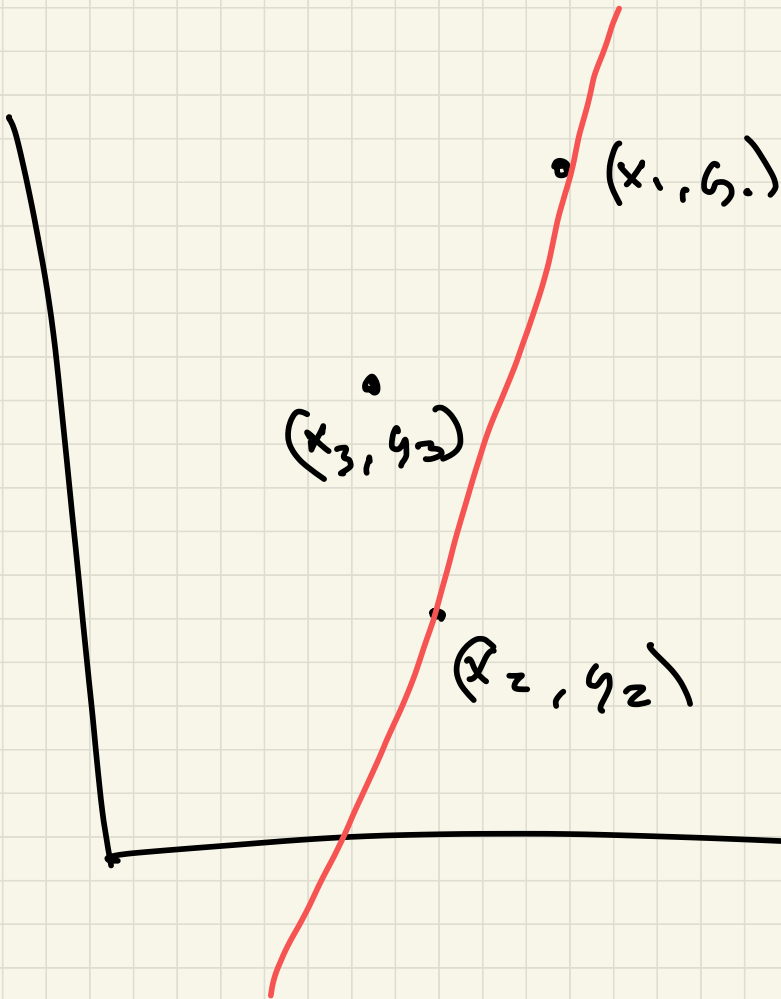
$$\alpha_0 = y_1 - \alpha_1 x_1$$

$$f_2(\alpha) = (y_2 - \alpha_0 - \alpha_1 x_2)^2$$

$$\alpha_0 = y_2 - \alpha_1 x_2$$

$$\alpha_0 = y - \alpha_1 x$$

$$\alpha = (\alpha_0, \alpha_1)$$



if $n \geq \#$ parameters
(not non-general case)

then f strongly
convex

$$f = \sum_{i=1}^n f_i$$

Batch Gradient Descent (Batch size = n)

$$\alpha = \alpha - \gamma \nabla f(\alpha)$$

$$f = \sum_{i=1}^n f_i$$

$$\nabla f(\alpha) = \nabla \sum_{i=1}^n f_i(\alpha) = \sum_{i=1}^n \nabla f_i(\alpha)$$

$$\frac{\partial}{\partial \alpha_j} f(\alpha) = \sum_{i=1}^n \frac{\partial}{\partial \alpha_j} f_i(\alpha) = 2 \sum_{i=1}^n (M_{\alpha}(x_i) - y_i) x_i^j$$

$$\nabla f(\alpha) = 2 \sum_{i=1}^n \underbrace{(M_{\alpha}(x_i) - y_i)}_{\text{residual}} \underbrace{\begin{pmatrix} 1 \\ x_i \\ x_i^2 \end{pmatrix}}_{\text{feature vector}}$$

Incremental Gradient Descent

$$\nabla f(x) \approx \nabla f_i(x) = 2(Ma(x_i) - y_i) (1, x_i, x_i^2)$$

$$\alpha^{(0)} = \alpha^{\text{start}} \quad ; \quad i = 1$$

repeat

$$\alpha^{(k+i)} = \alpha^{(k)} - \delta_k \nabla f_i(\alpha^{(k)})$$

$\mathcal{O}(1)$ time
to compute

$$i = (i+1) \bmod n$$

$$\text{until } (\|\nabla f_i(\alpha^{(k)})\| \leq \tau)$$

return $\alpha^{(k)}$

take average over B steps
 $B = 10?$

Stochastic Gradient Descent (SGD)

init: $\alpha^{(0)} = \alpha^{\text{start}}$

repeat

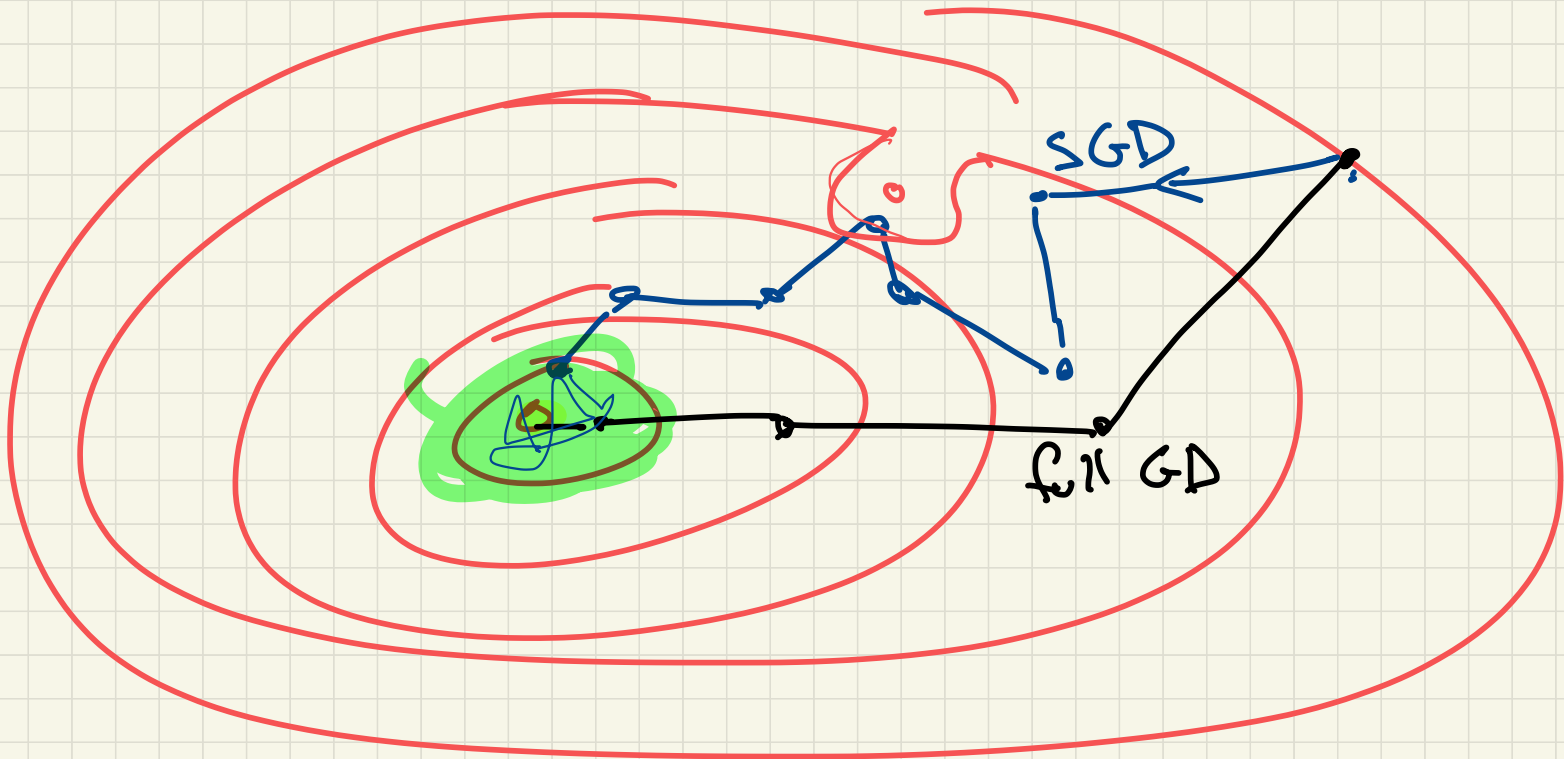
choose $i \in [1 \dots n]$ at random

$$\alpha^{(k+1)} = \alpha^{(k)} - \eta \nabla f_i(\alpha^{(k)})$$

until $(\|\nabla f(\alpha^{(k)})\| \leq \tau)$

\uparrow average

return $\alpha^{(k)}$



SGD : more steps
smaller steps | but each step is
much, much faster