

# L16: Stochastic GD

---

Applied on Data

Mar 4, 2026

FoDA

Jeff M. Phillips



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

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

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  $L((X, y), M_\alpha)$

$$\begin{aligned} f(\alpha) &= L((X, y), M_\alpha) = \text{SSE}((X, y), M_\alpha) \\ &= \sum_{(x_i, y_i) \in (X, y)} (M_\alpha(x_i) - y_i)^2 \end{aligned}$$

Model  $M_\alpha$  for polynomial regression

$$d=1 \quad x_i \in \mathbb{R} \quad p=2$$

$$M_\alpha(x_i) = \langle \alpha, (1, x_i, x_i^2) \rangle = \alpha_0 + \alpha_1 x_i + \alpha_2 x_i^2$$

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

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

consider  
 $n=1$

$$\frac{\partial}{\partial \alpha_j} f(\alpha) = \frac{\partial}{\partial \alpha_j} (M_\alpha(x_1) - y_1)^2 = 2(M_\alpha(x_1) - y_1) \frac{\partial}{\partial \alpha_j} (M_\alpha(x_1) - y_1)$$

$$= 2(M_\alpha(x_1) - y_1) \frac{\partial}{\partial \alpha_j} \left( \sum_{i=0}^2 \alpha_i x_1^i - y_1 \right)$$

$$= 2(M_\alpha(x_1) - y_1) x_1^j$$

$j=0,1,2$

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

$$n=1$$

$$\nabla f(\alpha) = \left( \frac{\partial}{\partial \alpha_0} f(\alpha), \frac{\partial}{\partial \alpha_1} f(\alpha), \frac{\partial}{\partial \alpha_2} f(\alpha) \right)$$

$$= 2 \underbrace{(M_\alpha(x_i) - y_i)}_{\text{scalar error = residual at } (x_i, y_i)} \underbrace{(1, x_i, x_i^2)}_{\text{data point}}$$

LMS update rule, Widrow-Hoff update

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

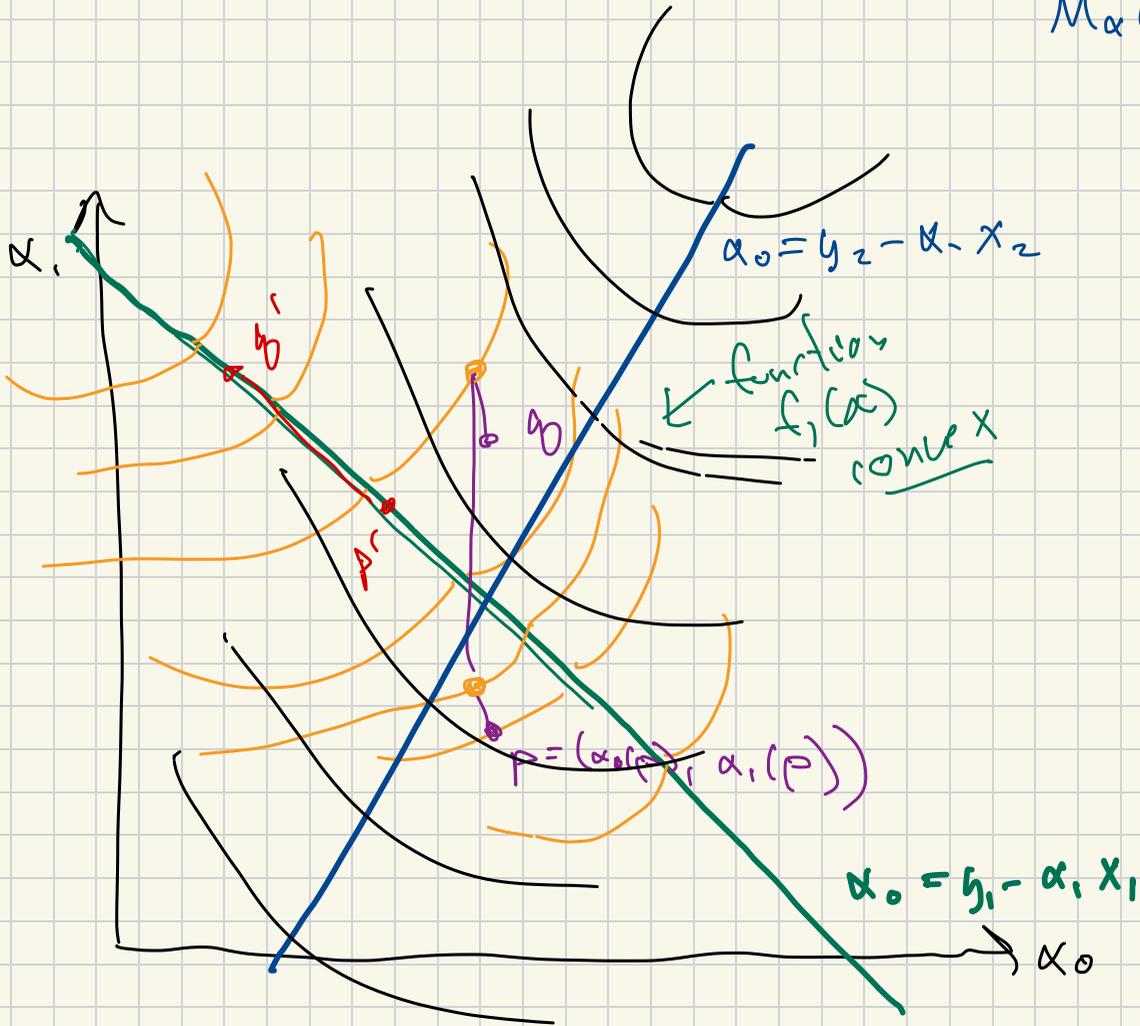
↖ for 1 data point  $(x_i, y_i)$

$n \geq 1$   
decomposable function

$$\begin{aligned} f_i(\alpha) &= (M_{\alpha}(x_i) - y_i)^2 \\ &= (\alpha_0 + \alpha_1 x_i + \alpha_2 x_i^2 - y_i)^2 \end{aligned}$$

if each  $f_i$  is convex, then

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



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

$$(x_1, y_1)$$

$$y_1 \approx \alpha_0 + \alpha_1 x_1$$

$$f(\alpha) = 0$$



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

$$(x_2, y_2)$$

$$f_2(\alpha) = 0 \implies \alpha_0 = y_2 - \alpha_1 x_2$$

$\hookrightarrow \exists f_i$  usually  
 strongly convex

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

$\alpha_0$

Batch

Gradient Descent

$n = \text{large } x$

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

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

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

$$\nabla f(\alpha) = 2 \sum_{i=1}^n (M_\alpha(x_i) - y_i) \begin{pmatrix} 1 \\ x_i \\ x_i^2 \end{pmatrix}$$

# Incremental Gradient Descent

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

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

repeat

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

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

until  $(\|\nabla f_i(\alpha^{(k)})\| < \epsilon)$

return  $\alpha^{(k)}$

average over  
e.g. 10 steps

# Stochastic Gradient Descent

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

repeat

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

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

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

return  $\alpha^{(k)}$

$\nwarrow$  average

