

FoDA

L17

- Gradient Descent
- Fitting Models to Data

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

Goal: Model M_α

$$\alpha = (\alpha_1, \dots, \alpha_k)$$

minimize loss function

$$f(\alpha) = L((X, y), M_\alpha) \quad f(\alpha): \mathbb{R}^k \rightarrow \mathbb{R}$$

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

$$(X, y) \quad X \in \mathbb{R}^n \quad y \in \mathbb{R}^n$$

$$(x_i, y_i) \in \mathbb{R} \times \mathbb{R}$$

model
quadratic

$$M_\alpha(x_i) = \langle \alpha, (1, x_i, x_i^2) \rangle$$

$$= \alpha_0 + \alpha_1 x_i + \alpha_2 x_i^2$$

Gradient Descent $f(\alpha) = \sum_{i=1}^n (M_\alpha(x_i) - y_i)^2$

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

Single Data Point $(x_i, y_i) \stackrel{2=1}{=} (x_i, g_i)$

$$f(\alpha) = f_1(\alpha) = (\alpha_0 + \alpha_1 x_i + \alpha_2 x_i^2 - g_i)^2$$

convex in α

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

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

$$= 2 (M_\alpha(x_i) - g_i) \frac{\partial}{\partial \alpha_j} (M_\alpha(x_i) - g_i)$$

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

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

LMS update rule / Widrow-Hoff learning rule

$$\nabla f(\alpha) = \left(\frac{\partial}{\partial \alpha_0} f(\alpha), \frac{\partial}{\partial \alpha_1} f(\alpha), \frac{\partial}{\partial \alpha_2} f(\alpha) \right) = 2 (M_\alpha(x_i) - g_i) (1, x_i, x_i^2)$$

Decomposable Functions

($n > 1$) Data Points

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

$$f_i(x) = (M_x(x_i) - y_i)^2$$

$$f(x) = \sum_{i=1}^n f_i(x) = SSE((x, y), M_x)$$

Batch Gradient Descent

$$\nabla f(x) = \sum_{i=1}^n \nabla f_i(x) = \left(\sum_{i=1}^n \frac{\partial}{\partial \alpha_0} f_i(x), \dots, \sum_{i=1}^n \frac{\partial}{\partial \alpha_k} f_i(x) \right)$$

LMS Update

$$= \sum_{i=1}^n z (M_x(x_i) - y_i) (1, x_i, x_i^2, \dots)$$

LMS Update

$$f(\alpha) = \sum_{i=1}^n (\alpha_0 + \alpha_1 x_i + \alpha_2 x_i^2 - y_i)^2$$
$$= \sum_{i=1}^n f_i(\alpha)$$

LMS Update

$$\nabla f_i(\alpha) = 2 (\alpha_0 + \alpha_1 x_i + \alpha_2 x_i^2 - y_i) \cdot (1, x_i, x_i^2)$$

$$\nabla f(\alpha) = \sum_{i=1}^n 2 (\alpha_0 + \alpha_1 x_i + \alpha_2 x_i^2 - y_i) (1, x_i, x_i^2)$$

residual $(\alpha_0 x_i - y_i)$

expanded explanatory data point $x_i \rightarrow (1, x_i, x_i^2)$

Since $f_i(\alpha)$ convex

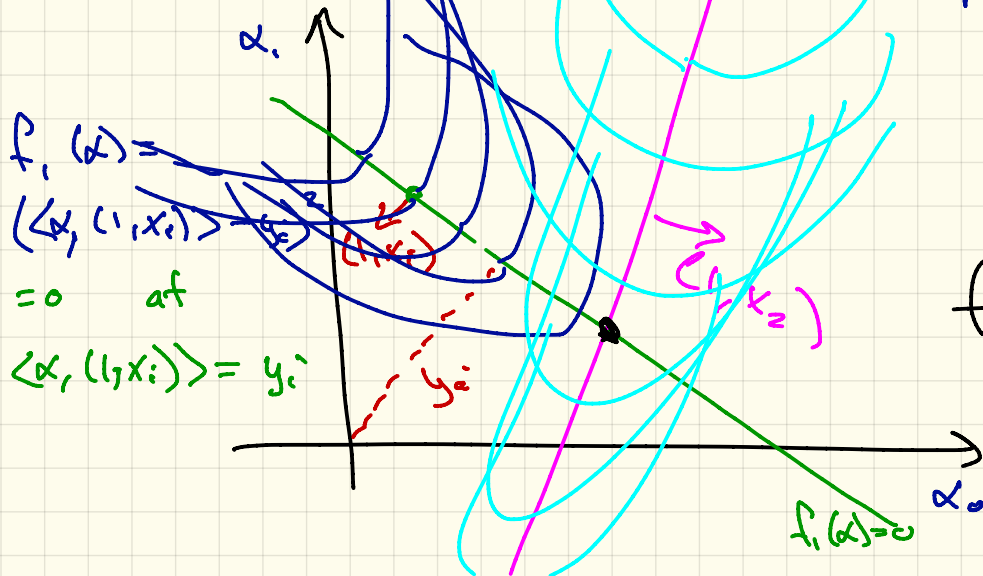
$\rightarrow \sum_{i=1}^m f_i(\alpha)$ convex

strongly convex?

linear regression

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

expand $x_i \rightarrow (1, x_i) \in \mathbb{R}^2$
 model $\alpha = (\alpha_0, \alpha_1)$



$f_1(\alpha) = \langle \alpha, (1, x_i) \rangle$
 $= 0$ at $\langle \alpha, (1, x_i) \rangle = y_i$

$$f(\alpha) = f_1(\alpha) + f_2(\alpha)$$

strongly convex

if $\begin{bmatrix} 1 & x_1 & x_1^2 \\ 1 & x_2 & x_2^2 \\ \vdots & \vdots & \vdots \end{bmatrix}$ full rank

What about Big Data? ^{n very big}

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

approximate $\nabla f(\alpha)$ in constant time

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

Incremental Gradient Descent (x_i, y_i)

0. Initialize $\alpha^{(0)} \in \mathbb{R}^d$ $i=1, k=0$

1. repeat

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

maybe make smaller than full GD
 $\nabla f_i(\alpha^{(k)})$ ← Grad at data point.

2. until $i = (i+1) \bmod n$
 $\|\nabla f(\alpha^{(k)})\| \leq T$

take sliding average over B step

Stochastic Gradient Descent (SGD)

0. $x^{(0)} = x \in \mathbb{R}^d$

1. repeat

a. Randomly choose $i \in \{1, 2, \dots, n\}$

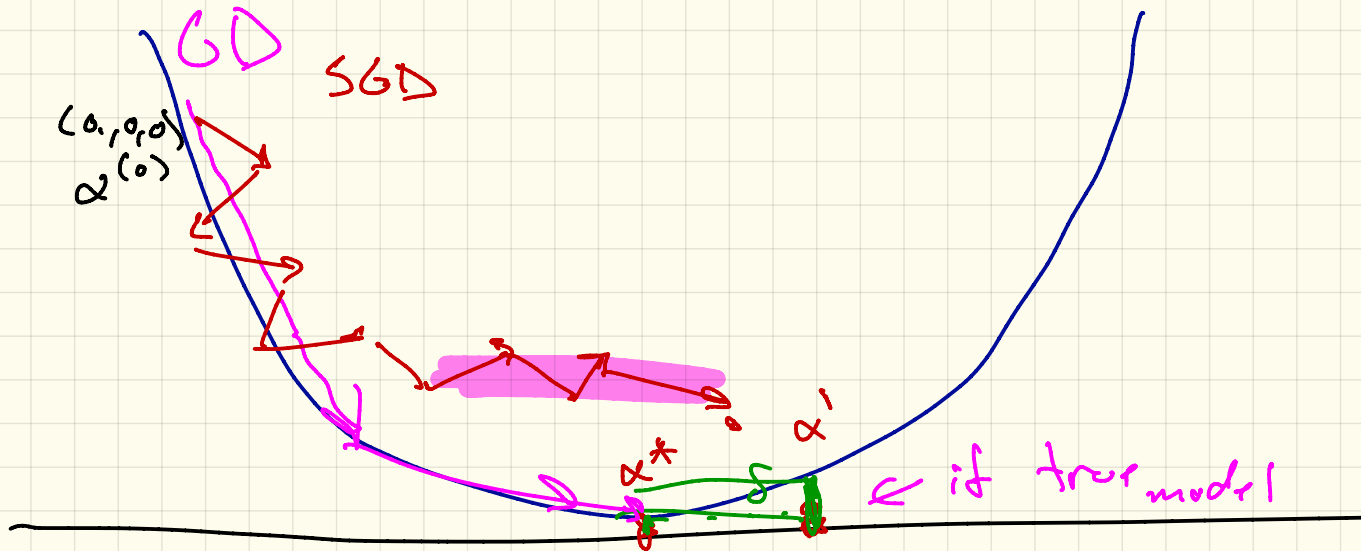
b. $x^{(k+1)} = x^{(k)} - \gamma \nabla f_i(x^{(k)})$

2. until $(\|\nabla f_i(x^{(k)})\| \leq \tau)$

3. Return $x^{(k)}$

... Recently

SGD tends to generalize better,
than full / Batch GD.



Strongly Convex

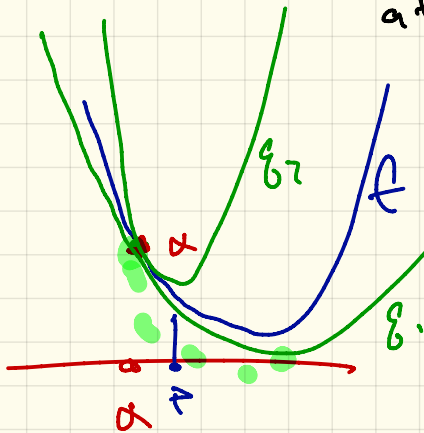
$f(x)$ is str. convex if

at each $x \in \mathbb{R}^d$ \exists quadratic
functions $g_1, g_2 : \mathbb{R}^d \rightarrow \mathbb{R}$

so all $p \in B_r(x)$ have
Grad of L-Lipschitz

$$g_1(p) \leq f(p) \leq g_2(p)$$

$$g_1(x) = f(x) = g_2(x)$$



$$g_1(p) = f(x) + \langle \nabla f(x), p - x \rangle + \frac{\mu}{2} \|p - x\|^2$$

Decomposable functions

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

usually each f_i is a simple form
1 dots point (x, y)

Non-decomposable f(x)

$$f(x) = (x_1 + 1 + x_2 x_3)^2 (x_1 - x_2)(x_3^2)$$

