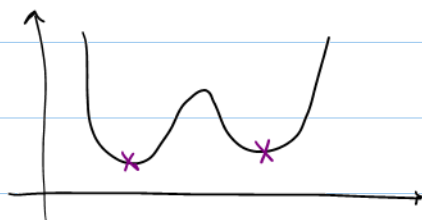


$$\vec{x} = (\alpha, \beta, \gamma, \dots)$$

مسئله optimization

$$\min [h(\vec{x})]$$

$$f(x) = x^4 - x^2$$



ناله

$$\vec{x} = \vec{x} + \alpha \vec{u}$$

$$\vec{u} = -\vec{\nabla} h(\vec{x})$$

جهت سیر برای سز درین مسرت:

(\*) حد نته هم:

$$\vec{x}_* = ?$$

①  $x$  اریه کی باسه = جواب هارا تغییر می دهه.

$$\vec{x}_{new} = \vec{x} - \alpha \vec{\nabla} h(\vec{x})$$

②  $\alpha$  کی ظاهر می شود محلا  $10^{-2}, -3$

③ tolerance  $\sim$  نزیب بول  $x$  و  $x_{new}$  ر نشان می دهه

④ iteration step  $n = \text{step}$  تهر تم ها تا min

بارسها هرحه می توانه باسه، مل  $x$  می توان وزن نودن ها باسه

# - Batch Gradient Descent (single or Mini batch).

برای مثال 2 داریم که با هم درجه و مرتبه regression

$$(x_i, y_i)$$

$x_1$	$y_1$
$x_2$	$y_2$
$\vdots$	$\vdots$
$x_m$	$y_m$

$$h(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m [y_i - (\theta_0 + \theta_1 x_i)]^2$$

number of training sample : m

input-output pairs :  $x_i, y_i$

Optimization parameters :  $\theta_0, \theta_1$

$$\begin{cases} \frac{\partial h}{\partial \theta_0} = -\frac{1}{m} \sum_{i=1}^m y_i - (\theta_0 + \theta_1 x_i) \\ \frac{\partial h}{\partial \theta_1} = -\frac{1}{m} \sum_{i=1}^m x_i [y_i - (\theta_0 + \theta_1 x_i)] \end{cases}$$

$$\theta_0 = \theta_0 - \alpha \frac{\partial h}{\partial \theta_0} \quad , \quad \theta_1 = \theta_1 - \alpha \frac{\partial h}{\partial \theta_1}$$

حال بقیه به چند نیت دارد که Run می‌کنیم

این از کل m ها استفاده کنیم = Batch Gradient Descent ①

اگر در update، در هر m جمع می‌کنیم = به این واسطه راحت‌تر می‌شود انتخاب کنیم  
② single batch gradient descent = انت‌دخیر دارد

چون انتخاب آسان بدنه Random است ← stochastic است

حالت این تعداد  $T$  تا برداریم = mini batch gradient descent (3)

معمولاً  $T$  را 32 یا 64 می‌گیرند ← parallel ← GPU می‌تواند بد

حالت چگونگی  $T$  تا انتخاب کنیم : chunk برداریم = بست سرهم

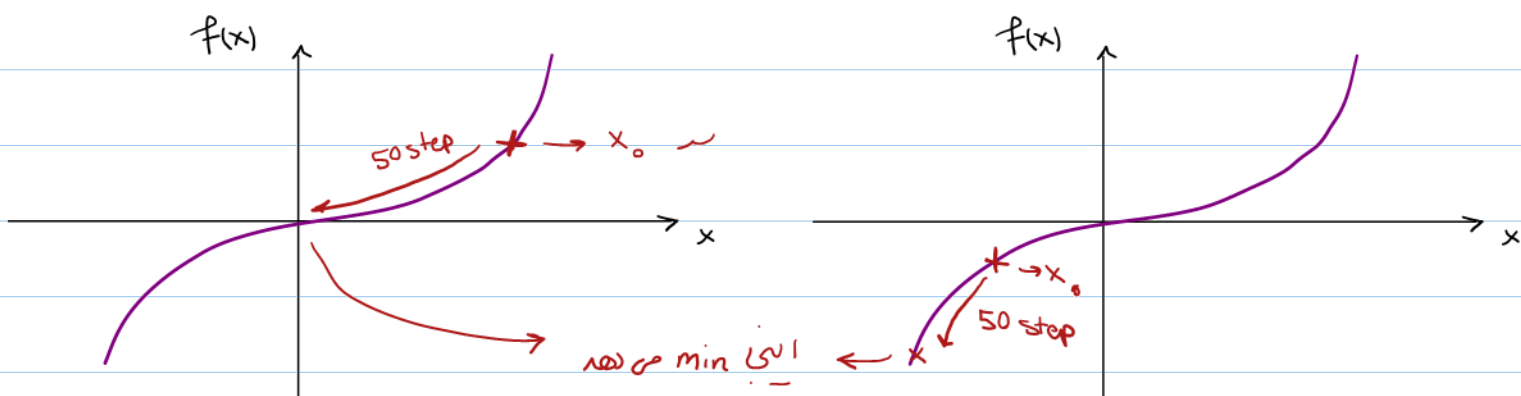
Stochastic Random برداریم =

مثال: تابع Convex یا غیر Convex چیست؟

اگر اطلاعات از تابع Convex باشد و اطلاعات را داریم = توان می‌سازیم پس می‌خواهد

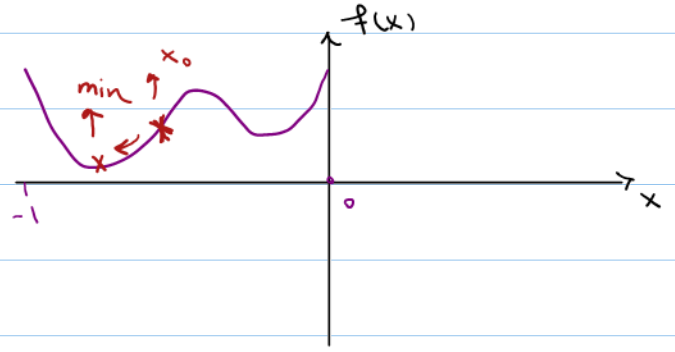
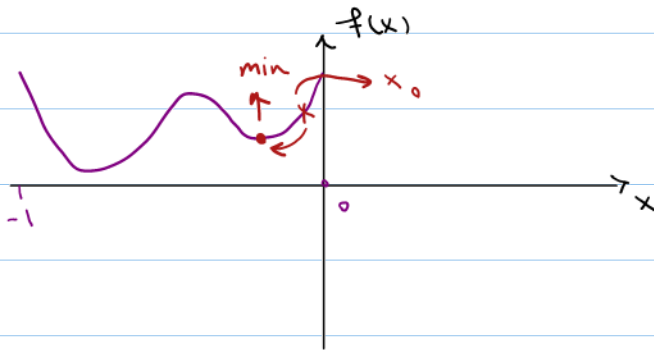
حالت این Non-Convex باشد چند حالت دارد

Saddle point:  $f(x) = x^3$   $f'(x) = 3x^2$  (مستقیم است اما min نیست)



این واضح است است و غیر تکرار

$$f(x) = x^3 + \sin(10x)$$



در این مسئله هدف ما پیدا کردن  $\min$  است. ما می‌خواهیم بدانیم که آیا این مسئله را می‌توانیم حل کنیم یا نه.

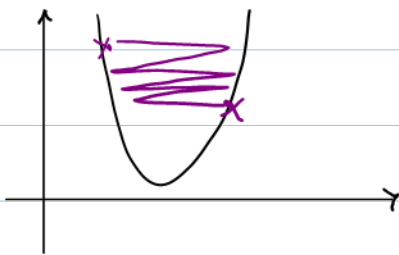
در این مسئله ما می‌خواهیم بدانیم که آیا این مسئله را می‌توانیم حل کنیم یا نه.

ما می‌خواهیم بدانیم که آیا این مسئله را می‌توانیم حل کنیم یا نه.

برای حل این مسئله ما می‌خواهیم بدانیم که آیا این مسئله را می‌توانیم حل کنیم یا نه.

ما می‌خواهیم بدانیم که آیا این مسئله را می‌توانیم حل کنیم یا نه.

ما می‌خواهیم بدانیم که آیا این مسئله را می‌توانیم حل کنیم یا نه.



ما می‌خواهیم بدانیم که آیا این مسئله را می‌توانیم حل کنیم یا نه.

Algorithmic variations of GD:

1) Momentum - Steering:

$$x = x - \alpha \underbrace{\nabla_{\vec{x}} L(x)}_{\vec{u}}$$

در این مسئله ما می‌خواهیم بدانیم که آیا این مسئله را می‌توانیم حل کنیم یا نه.

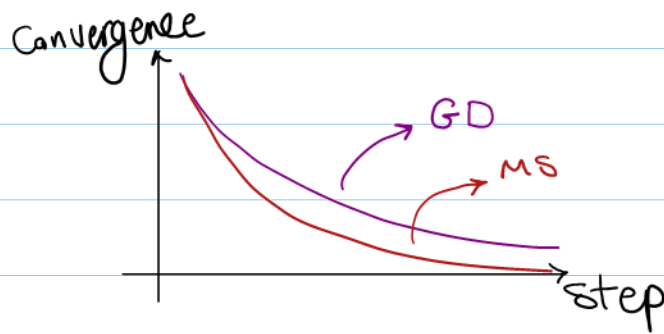
$$x_{s+1} = x_s + \alpha \vec{u}_s$$

حال میخواهیم اینجا به  $u_s$  از step قبلی اطلاعات بدیم

$$u_s = \gamma u_{s-1} - \alpha \nabla_x h(x)$$

$\gamma = 0.9$  typically

چون از قبلی هم میبرد نامش momentum است. علاوه بر سرعت convergence، این را هم کاهش میدهد.



## 2) Nesterov Accelerated Gradient (NAG):

$$u_s = \gamma u_{s-1} - \alpha \nabla_x h(x - \gamma u_{s-1})$$

convergence سریعتره.

## 3) Annealing & learning rate schedules

هنگامی داریم که convergence خیلی کند میسریم،  $\alpha$  را کاهش می دهیم.

اصطلاحاً میگویند سرد کردن ← cooling

چند روش برای کاهش وجود دارد.

### 3-1) Ada-Grad (Adaptive Gradient Alg).

$$\frac{\partial h}{\partial x_i} \rightarrow g_i^s \quad g_i^{(0)}, \dots, g_i^{(s-1)}$$

$$G^s = (g_i^0)^2 + \dots + (g_i^{s-1})^2$$

$$x_i = x_i - \frac{\alpha}{\sqrt{G^s + \epsilon}} g_i^s$$

↓ که اینر  $G$  صغیرتره مشکل ای داره

برخ صورتی که  $\epsilon$  step قبل رانده داریم نه کمه را  $\epsilon$  بی این است که در  $\epsilon$  step نزدیک Converge  $G$  ها خیلی بزرگ صورتی

### 3-2) RMSprop root mean square propagator:

$$g_s = (\nabla h^2)^{(s)}$$

$$g_s^2 = \beta g_{s-1}^2 + (1 - \beta) g_s^2$$

0.9 = typically

### 3-3) Ada Delta : an adaptive learning rate method

$$x_i^{s+1} = x_i^s - \frac{\text{RMS}[\Delta x_i]^{s-1}}{\text{RMS}[g_i]^s} g_i^{(s)}$$

این  $\Delta x$  ها بزرگ هسته  $\leftarrow$  rate ها بزرگ هسته و کم با بخش می یابید

#### 4) Adam: adaptive moment estimation

$$g = \nabla_x l$$

putting  $\nabla$  as  $z \Rightarrow m = \langle z \rangle$ ,  $v = \langle z^2 \rangle$

$$\begin{cases} m_s = \beta_1 m_{s-1} + (1 - \beta_1) g_s & \beta_1 = 0.9 \\ v_s = \beta_2 v_{s-1} + (1 - \beta_2) g_s^2 & \beta_2 = 0.999 \end{cases}$$

$$\Rightarrow \begin{cases} \hat{m}_s = \frac{m_s}{1 - \beta_1^t} \\ \hat{v}_s = \frac{v_s}{1 - \beta_2^t} \end{cases} \Rightarrow \text{Bias corrector}$$

$$x^{s+1} = x^s - \frac{\alpha}{\sqrt{\hat{v}_s + \epsilon}} m^s$$

$$m_0, v_0 = 0, \quad \alpha \approx 10^{-3}, \quad \epsilon \approx 10^{-8}$$

Exercise 1:

من خواهم جواب حساب کنم

$$Ax = B$$

تمرین اول

داده ماتریس  $A$  و بردار  $B$

من خواهم با optimization حل کنم . Norm چیست ؟

$$h(x) = \min \|Ax - B\|$$

آه بنویسم نه با GD حل کنم این و \*

در جایی امکان دارد  $full\ rank$  نباشد، یعنی یک شرط ضعیف دیگری است.

در این  $\Rightarrow$   $inverse\ det \rightarrow \infty$  می شود

اگر در  $inverse$  ماتریس  $A$  دچار مشکل بشویم تمام الگوریتم ها از کار می افتد.

5 راهی که این مشکل را حل می کنند باید  $Optimization$  را برای یکس انجام دهید.

توضیح دهید  $\leftarrow$  برای  $det \approx 10^{-4}$

حال برای حل مسائلی که به نیتیم چند مرحله وجود دارد.

- به صورت  $uniform$  در بازه  $?$  اولیه ریزیم بهم  $\Rightarrow$

- به  $u$  جهت ریزیم بهم  $\Rightarrow$  سانس داده برپرد برخط بردار  $\Rightarrow$

نات  $n = steps = epochs$

$\alpha = learning\ rate = 0.005$

momentum  $\gamma = 0.9$

Noise intensity  $y(\sigma) \Rightarrow u = -\nabla h(x) \rightarrow -\nabla h(x) + \sigma \epsilon(0,1)$   
white noise  $\epsilon$

Noise decay  $= 0.99$   $\sigma^{s+1} = 0.99 \sigma^s$

حال مقادیر را بهترین ها را save می کنیم



کدگذاری: برای fit کردن معادله خط  $\theta_0, \theta_1$  (code SGD2)

کد دوم: برای این تابع با هر شرایطی بارش نسبت به  $\min$  درست را می‌دهد. (code SGD3)

برای تست کردن optimization برای چند بعد، می‌توانیم از تست زیر استفاده کنیم.

چون  $\min$  در دره بسیار ریز است. Rosenbrock function

و جواب  $\min$  و saddle آن را بدست می‌دهیم.

می‌توان برای fine tune کردن پارامترها استفاده کرد.