



بهینه سازی  
مبانی بهینه سازی نامقید  
روش گرادیان نزولی، روش نیوتن

محسن هوشمند  
دانشکده تکنولوژی اطلاعات و علم رایانه  
دانشگاه تحصیلات تکمیلی علوم پایه زنجان

# مقدمه

یافتن کمینه تابع هدف

مبتنی بر متغیرهای حقیقی بدون قید

از تمامی فضای اعداد حقیقی

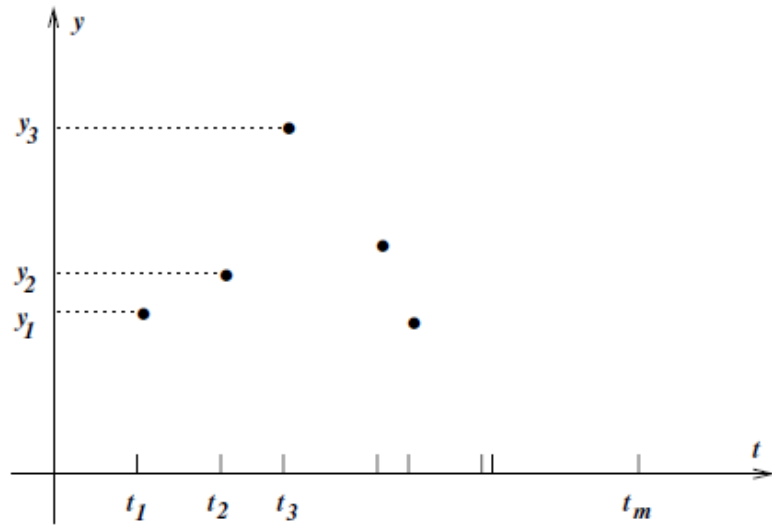
$$\min_x f(x)$$

$$x \in \mathbb{R}^n$$

$$n \geq 1$$

$$f: \mathbb{R}^n \rightarrow \mathbb{R} \text{ تابعی هموار}$$

# مقدمه - مثال



یافتن منحنی که چند داده را تقریب بزند

$$\phi(t; x) = x_1 + x_2 e^{-(x_3 - t)^2 / x_4} + x_5 \cos(x_6 t) \quad \text{نمایش نمایی و نوسانی}$$

$x_i$  -ها ضرائب مدل

▪ مجهول‌ها

$$x = (x_1, x_2, \dots, x_6)^T$$

تعریف باقی‌مانده

$$r_j(x) = y_j - \phi(t_j; x), \quad j = 1, 2, \dots, m$$

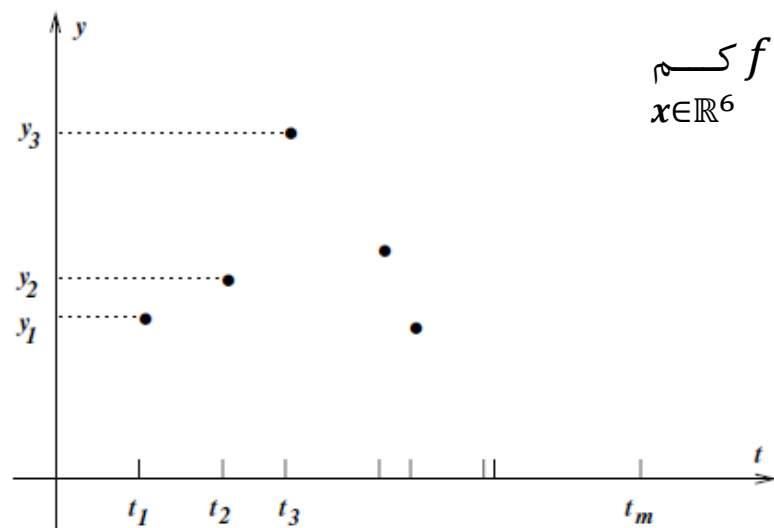
▪ نمایشگر تفاوت بین مدل و داده مشاهده شده

$$f(x) = r_1^2(x) + r_2^2(x) + \dots + r_m^2(x) \quad \text{کم}$$

$x \in \mathbb{R}^6$

▪ مسئله برازش کمترین مربعات غیرخطی

# مقدمه - مثال - ادامه



$$\min_{x \in \mathbb{R}^6} f(x) = r_1^2(x) + r_2^2(x) + \dots + r_m^2(x)$$

- مسئله برازش کمترین مربعات غیرخطی
- از انواع بهینه‌سازی نامقید

$$x^* = (1.1, 0.01, 1.2, 1.5, 2.0, 1.5)$$

فرض: پاسخ

$$f(x^*) = 0.34$$

مقدار کمینه

وجود تفاوت بین مقدار واقعی و مقدار بدست‌آمده

▪ چگونه  $x^*$  را کمینه‌ساز  $f$  می‌دانیم

▪ نیاز به تعریف «راه‌حل»

# تشخیص کمینه محلی

کمینه‌ساز سراسری  $f$

- نقطه بدست‌دهنده کمترین مقدار تابع  $f(x^*) \leq f(x) \forall x$
- مشکل در یافتن
- اطلاع صرفاً از دور و بر هر نقطه و نه بیشتر
- عدم اطمینان از وجود شکاف عمیق در داده
- بیشتر الگوریتم‌ها یابنده کمینه‌ساز محلی

کمینه‌ساز محلی  $x^*$

- اگر در همسایگی  $x^*$  (نمایش با  $\mathcal{H}$ )  $f(x^*) \leq f(x), x \in \mathcal{H}$
- معروف به کمینه‌ساز محلی ضعیف
- کمینه‌ساز محلی اکید
- $f(x^*) < f(x), x \in \mathcal{H}, x \neq x^*$

# تشخیص کمینه محلی

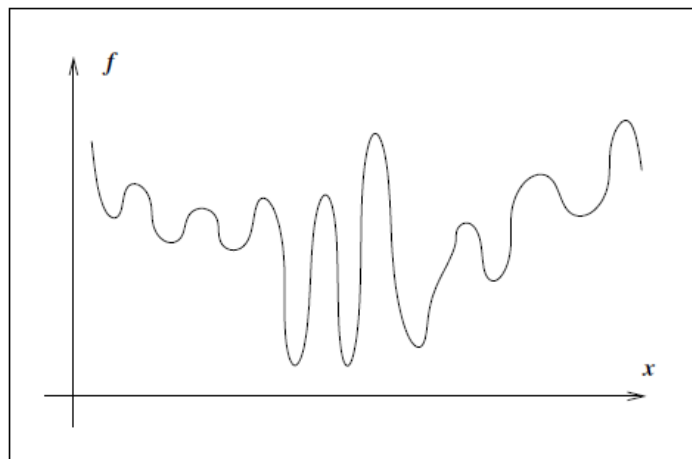
مثال

$$f(x) = 2 \cdot$$

مثال ۲

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

کمینه‌ساز محلی تک‌افتاده



# یافتن کمینه محلی

جستجو همسایگی؟

تابع هموار

▪ روش‌های کارا تر و عملی تر

تابع مشتق‌پذیر مرتبه اول و دوم

▪ وجود گرادیان و هسی

▪ محتملا بتوان نقطه‌ای را کمینه محلی (کمینه محلی قوی) نامید

▪ ابزار ریاضی مطالعه کمینه‌سازها

▪ قضیه تیلور

قضیه تیلور

# قضیه تیلور

$f: \mathbb{R}^n \rightarrow \mathbb{R}$  تابعی مشتق پذیر پیوسته و  $\mathbf{p} \in \mathbb{R}^n$ ، آن گاه

$$f(\mathbf{x} + \mathbf{p}) = f(\mathbf{x}) + \nabla f(\mathbf{x} + t\mathbf{p})^T \mathbf{p}$$

$$t \in (0,1)$$

همچنین، اگر مشتق پذیر مرتبه دوم پیوسته باشد، آن گاه

$$\nabla f(\mathbf{x} + \mathbf{p}) = \nabla f(\mathbf{x}) + \int_0^1 \nabla^2 f(\mathbf{x} + t\mathbf{p}) \mathbf{p} dt$$

$$f(\mathbf{x} + \mathbf{p}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T \nabla^2 f(\mathbf{x} + t\mathbf{p}) \mathbf{p}$$

$$t \in (0,1)$$



# قضیه شروط لازم مرتبه اول

$\nabla f(\mathbf{x}^*) = 0$  آن گاه  $\mathbf{x}^*$  کمینه محلی و  $f$  در همسایگی نقطه مذکور مشتق پذیر و پیوسته است، اثبات

$\nabla f(\mathbf{x}^*) = 0$  اگر  $\mathbf{x}^*$  نقطه ماناست نتیجه قضیه: هر کمینه محلی، نقطه ماناست.

# قضیه شروط لازم مرتبه دوم

$x^*$  کمینه محلی  $f$  و  $\nabla^2 f$  موجود و پیوسته همسایه نقطه مذکور، آن گاه  $\nabla f(x^*) = 0$  و  $\nabla^2 f(x^*)$  مثبت نیمه معین است.

اثبات

# قضیه شروط کافی مرتبه دوم

$\nabla^2 f$  موجود و در همسایگی  $x^*$  پیوسته و  $\nabla f(x^*) = 0$  و  $\nabla^2 f(x^*)$  مثبت معین، آن گاه نقطه مذکور کمینه محلی اکید تابع  $f$  است.

اثبات

# شرایط بهینگی کمینه‌سازی چندمتغیره

قضیهٔ شروط لازم کمینه محلی ضعیف

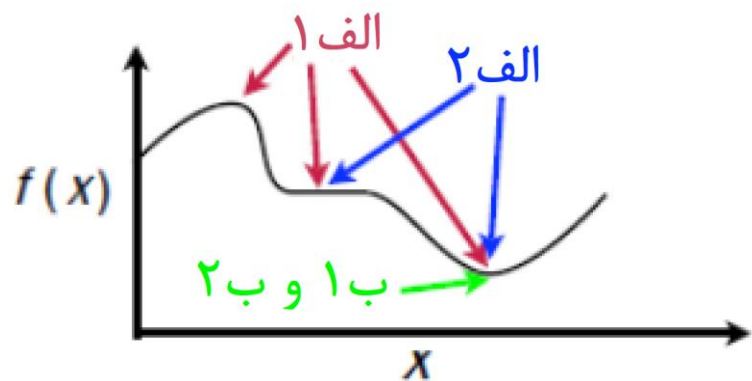
الف ۱:  $\nabla f(x^*) = 0$  نقطه مانا

الف ۲:  $\nabla^2 f(x^*)$  مثبت نیمه‌معین

قضیهٔ شروط کافی کمینه محلی قوی

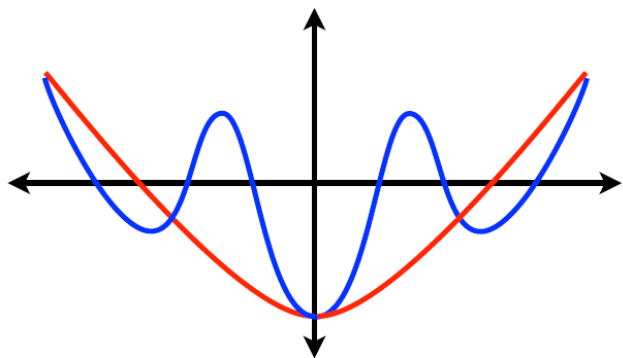
ب ۱:  $\nabla f(x^*) = 0$

ب ۲:  $\nabla^2 f(x^*)$  مثبت معین



# قضیه کوژ

$f$  کوژ باشد، هر  $x^*$  کمینه‌ساز محلی کمینه‌ساز سراسری  $f$  است. همچنین اگر  $f$  مشتق‌پذیر باشد، آن‌گاه هر نقطه مانای  $x^*$  کمینه‌ساز سراسری  $f$  خواهد بود.



کوژی

▪ ماتریس هسی مثبت نمیه‌معین

اکیدا کوژ

▪ ماتریس هسی مثبت معین

# حل

نتایج بدست آمده از حسابان ساده و مقدماتی  
بنیادی جهت الگوریتم‌های بهینه‌سازی نامقید  
هر الگوریتم با یکی از روش‌ها به دنبال یافتن نقطه‌ای که گرادیان  $f$  ناپدید می‌شود

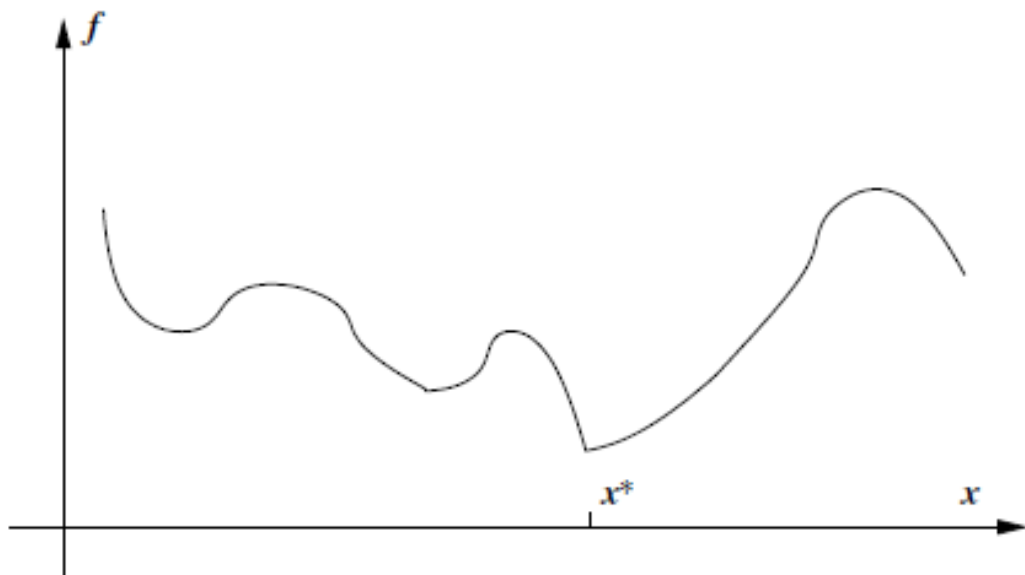
# مسائل ناهموار

توابع ناهموار و ناپیوسته

روشی عمومی برای حل وجود ندارد

در صورت اتصال چند قطعهٔ هموار با ناپیوستگی بین قطعه‌ها  
▪ امکان یافتن کمینه‌ساز با کمینه‌کردن جداگانهٔ هر قطعه

از روش‌های زیرگرادیان یا گرادیان تعمیمی



# الگوریتم‌های بهینه‌سازی نامقید

شصت سال اخیر

جملگی با شروع از نقطهٔ آغاز

تکرار دنباله‌ای از مراحل

پایان

- بهبود غیرممکن

- رسیدن به تخمین مناسبی از پاسخ

نزولی  $f(x_k) < f(x_{k-m})$

دو استراتژی اساسی جهت حرکت از نقطه فعلی  $x_k$  به نقطهٔ  $x_{k+1}$

- جستجو خط

- منطقه اعتماد

انتخاب جهت حرکت و اندازهٔ حرکت



# الگوریتم‌های بهینه‌سازی نامقید

روش‌های گرادیان-محور  
▪ انتخاب جهت حرکت و اندازه حرکت

الگوریتم گرادیان-محور

انتخاب مقدار اولیه  $x_0$  و  $k = 0$   
تازمان همگرا نشدن  
}

انتخاب جهت  $p_k$  و اندازه قدم  $\alpha_k$   
$$x_{k+1} = x_k + \alpha_k p_k$$
$$k \leftarrow k + 1$$
  
{

# دو استراتژی - جستجو خط و منطقه اعتماد

## استراتژی جستجو خط

- انتخاب جهت  $\mathbf{p}_k$
- جستجو در راستای جهت یافت شده از  $\mathbf{x}_k$  فعلی به تکرار جدیدی با مقدار  $f$  کمتر
- با حل مسئله کمینه‌سازی یک‌بعدی زیر جهت یافتن  $\alpha$

$$\min_{\alpha > 0} f(\mathbf{x}_k + \alpha \mathbf{p}_k)$$

## استراتژی منطقه اعتماد

- تعریف شعاعی  $\Delta > 0$  اطراف  $\mathbf{x}_k$  به عنوان منطقه اعتماد
- جمع‌آوری اطلاعات جهت ایجاد مدلی (تخمینی) از تابع  $f$
- تابع تخمین با نام  $m_k$
- دارای رفتار مشابه در نزدیکی نقطه  $\mathbf{x}_k$

$$\min_p m_k(\mathbf{x}_k + \mathbf{p})$$

$\mathbf{x}_k + \mathbf{p}$  در داخل منطقه اعتماد

$$\|\mathbf{p}\|_2 \leq \Delta$$

# دو استراتژی - ادامه

استراتژی منطقه اعتماد

▪ تابع تخمین  $m_k$  معمولا برابر با

$$m_k(\mathbf{x}_k + \mathbf{p}) = f_k + \mathbf{p}^T \nabla f_k + \mathbf{p}^T B_k \mathbf{p}$$

$B_k$  یا تابع هسی  $\nabla^2 f_k$  یا تخمینی از آن

# دو استراتژی - جستجو خط

## انتخاب $p_k$ در جهت کاهش

- بیشترین نزول (گرادیان نزولی): مرتبه اول، همگرایی خطی
- روش گرادیان مزدوج: مرتبه اول، همگرایی (سریعتر) خطی
- روش نیوتن: مرتبه دوم، همگرایی درجه دو
- روش شبه نیوتن: مرتبه اول تا تخمین مرتبه دوم، همگرایی ابرخطی
- با حل مسئله کمینه‌سازی یک‌بعدی زیر جهت یافتن  $\alpha$

## انتخاب طول قدم $\alpha_k$ برآورده‌گر شروط وولف

- کروسه‌گذاری: یافتن بازه‌ای شامل طول قدم مناسب
- نیمه‌سازی/درون‌یابی: محاسبه قدم مناسب در بازه فعلی

# شدیدترین نزول

شبیه پائین آمدن از کوه

انتخاب سریعترین جهت سرازیری

$$p_k = -\nabla f$$

کاربرد

- بهینه‌سازی
- شبکه عصبی
- یادگیری ماشین
- شبکه عمیق

مزایا

- صرفا استفاده از اطلاع مرتبه اول
- همیشه جهت نزول
- حافظه کم

معایب

- کند در مسائل سخت
- حساس به مقیاس گذاری

# شدیدترین نزول - مثال

حل تحلیلی

$$f(x) = x^2 - 2x + 2 = (x - 1)^2 + 1$$

$$\frac{df}{dx}$$

$$= 2x - 2$$

$$x = 1$$

# شدیدترین نزول - مثال - ادامه

حل با استفاده از شدیدترین نزول

$$f(x) = x^2 - 2x + 2 = (x - 1)^2 + 1$$

$$\frac{df}{dx}$$

مقدار بهینه را نمی‌دانیم  
▪ شروع از مقدار تصادفی  $x=3$

قدم اول: مشتق‌گیری

$$\frac{df}{dx} = 2x - 2$$

# شدیدترین نزول - مثال - /د/مه

قدم دوم

▪ مطالعه مشتق در نقطه داده شده

$$f'(3) = 2(3) - 2 = 4$$

▪ مشتق در کمینه باید صفر باشد

▪ مقدار مثبت مشتق

▪ نمایشگر اینکه مقدار تابع افزایشی است. و باید عقب رفت

▪ اگر مقدار  $x = -1$  به عنوان حدس اولیه انتخاب می‌شد

▪ آن‌گاه مشتق  $f' = -4$

▪ نمایشگر نزولی بودن مقدار و در نتیجه نیاز به جلو رفتن

مقدار فعلی مشتق نشانگر مسیر

▪ نزدیک شدن به یا دور شدن از کمینه



# شدیدترین نزول - مثال - ادامه

$$x_{i+1} = x_i + \alpha(-f'(x_i))$$

$x_{i+1}$  حدس بعدی

طول قدم  $\alpha = 0.2$

$$x_0 = 3$$

$$x_{i+1} = x_i + 0.2(-f'(x_i))$$

$$x_1 = 3 - \alpha f'(3) = 3 + 0.2(-4) = 2.2$$

$$x_2 = 2.2 - \alpha f'(2.2) = 2.2 + 0.2(-2.4) = 1.72$$

ادامه محتملا به پاسخ می‌رسد.

▪ ?

چرا شدیدترین نزول به جای روش تحلیلی و حسابان

نوشتن برنامه

## شدیدترین نزول - مثال ۲

یافتن بهترین خط برازش

$$y = \alpha x + \beta$$

x	y
1	1
2	1
2	2
3	2

به دنبال یافتن  $\alpha$  و  $\beta$

ابتدا نیاز به تعریف خطا

▪ خطا تفاضل بین داده و مقدار بدست آمده از مدل

$$e_1 = \hat{y}_1 - y_1$$

$$e_2 = \hat{y}_2 - y_2$$

نیاز به کل خطا

▪ یک روش: جمع تمامی مقادیر خطا  $e_{\text{کل}} = e_1 + e_2 + e_3 + e_4$

▪ غلطانداز

▪ روش بهتر  $e_{\text{کل}} = e_1^2 + e_2^2 + e_3^2 + e_4^2$

## شدیدترین نزول - مثال ۲ - ادامه

x	y
1	1
2	1
2	2
3	2

روش بهتر  $e_{\text{کل}} = e_1^2 + e_2^2 + e_3^2 + e_4^2$

$$e_1 = \hat{y}_1 - y_1 \Rightarrow e_1 = [\alpha x_1 + \beta] - y_1$$

$$e_2 = \hat{y}_2 - y_2 \Rightarrow e_2 = [\alpha x_1 + \beta] - y_2$$

$$e_3 = \hat{y}_3 - y_3 \Rightarrow e_3 = [\alpha x_1 + \beta] - y_3$$

$$e_4 = \hat{y}_4 - y_4 \Rightarrow e_4 = [\alpha x_1 + \beta] - y_4$$

خطای کل

$$\begin{aligned} e_{\text{کل}} &= \sum_{i=1}^4 e_i^2 \\ &= \sum_{i=1}^4 ([\alpha x_i + \beta] - y_i)^2 \end{aligned}$$

# شدیدترین نزول - مثال ۲ - ادامه

به دنبال یافتن خطی با کمترین خطای ممکن

x	y
1	1
2	1
2	2
3	2

$$e_{\text{کل}} = \sum_{i=1}^4 e_i^2 = \sum_{i=1}^4 ([\alpha x_i + \beta] - y_i)^2$$
$$\underset{\alpha, \beta}{\text{م}} \sum_{i=1}^4 ([\alpha x_i + \beta] - y_i)^2$$

استفاده از شدیدترین نزول

$$\begin{bmatrix} \alpha_{\text{ب}} \\ \beta_{\text{ب}} \end{bmatrix} = \begin{bmatrix} \alpha_{\text{ق}} \\ \beta_{\text{ق}} \end{bmatrix} - 0.2 \begin{bmatrix} \frac{\partial f}{\partial \alpha}(\alpha_{\text{ق}}, \beta_{\text{ق}}) \\ \frac{\partial f}{\partial \beta}(\alpha_{\text{ق}}, \beta_{\text{ق}}) \end{bmatrix}$$

# شدیدترین نزول - مثال ۲ - ادامه

x	y
1	1
2	1
2	2
3	2

$$f(\alpha, \beta) = \sum_{i=1}^4 ([\alpha x_i + \beta] - y_i)^2$$

$$\frac{\partial f}{\partial \alpha} = \sum_{i=1}^4 2([\alpha x_i + \beta] - y_i)x_i$$

$$\frac{\partial f}{\partial \beta} = \sum_{i=1}^4 2([\alpha x_i + \beta] - y_i)$$

$$\begin{bmatrix} \alpha_0 \\ \beta_0 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

مقدار اولیه

## شدیدترین نزول - مثال ۲ - ادامه

x	y
1	1
2	1
2	2
3	2

$$\begin{bmatrix} \alpha_{\text{ب}} \\ \beta_{\text{ب}} \end{bmatrix} = \begin{bmatrix} \alpha_{\text{ق}} \\ \beta_{\text{ق}} \end{bmatrix} - 0.2 \begin{bmatrix} \sum_{i=1}^4 2([\alpha x_i + \beta] - y_i)x_i \\ \sum_{i=1}^4 2([\alpha x_i + \beta] - y_i) \end{bmatrix}$$

$$\begin{bmatrix} \alpha_0 \\ \beta_0 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

مقدار اولیه

# شدیدترین نزول - مثال ۲ - ادامه

x	y
1	1
2	1
2	2
3	2

تحریر محل نزاع

- چهار داده
- محاسبات فراوان
- حال اگر یک میلیون داده موجود باشد
- یک میلیارد داده چه؟
- راه حل؟

$$\begin{bmatrix} \alpha_{\text{ب}} \\ \beta_{\text{ب}} \end{bmatrix} = \begin{bmatrix} \alpha_{\text{ق}} \\ \beta_{\text{ق}} \end{bmatrix} - 0.2 \begin{bmatrix} \sum_{i=1}^4 2([\alpha x_i + \beta] - y_i)x_i \\ \sum_{i=1}^4 2([\alpha x_i + \beta] - y_i) \end{bmatrix}$$

# شدیدترین نزول - مثال ۲ - ادامه

راه حل

▪ گرادیان نزولی تصادفی

$$\begin{bmatrix} \alpha_{j+1} \\ \beta_{j+1} \end{bmatrix} = \begin{bmatrix} \alpha_j \\ \beta_j \end{bmatrix} - 0.2 \begin{bmatrix} 2([\alpha_j x_1 + \beta_j] - y_1)x_1 \\ 2([\alpha_j x_1 + \beta_j] - y_1) \end{bmatrix}$$

$$\begin{bmatrix} \alpha_{j+2} \\ \beta_{j+2} \end{bmatrix} = \begin{bmatrix} \alpha_{j+1} \\ \beta_{j+1} \end{bmatrix} - 0.2 \begin{bmatrix} 2([\alpha_{j+1} x_2 + \beta_{j+1}] - y_2)x_2 \\ 2([\alpha_{j+1} x_2 + \beta_{j+1}] - y_2) \end{bmatrix}$$

- در هم کردن داده‌ها و خواندن از نمونه نخست
- ادامه تا خواندن همه نمونه‌ها و شروع کار با نمونه اول
- ادامه تا همگرایی



# شدیدترین نزول - مثال ۲ - ا/د/مه

## گرادیان نزولی

- کندتر
- صحیح‌تر
- دسته‌ای
- استفاده از همه نمونه‌ها

## گرادیان نزولی تصادفی

- سریع‌تر
- دقت کمتر
- استفاده از یک نمونه در هر زمان

## گرادیان نزولی زیردسته‌ای

- بیشتر از یک نمونه در هر زمان و کمتر از تمامی نمونه‌ها

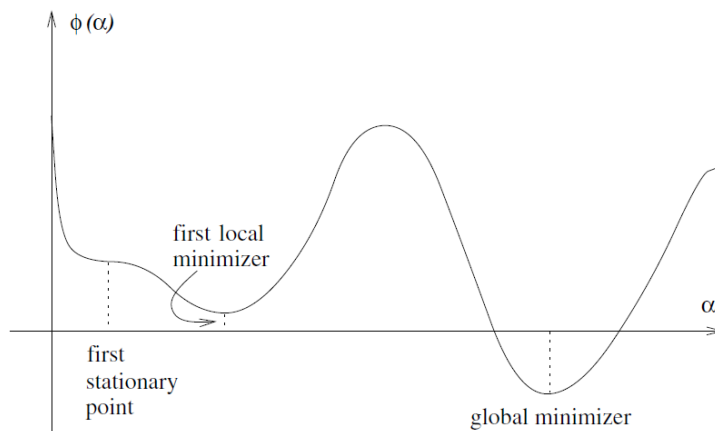
# طول قدم $\alpha$

## طول قدم

- نیاز به سبک‌سنگین کردن
- انتخاب طولی با کاهش مقدار قابل توجه  $f$
- بدون سپری کردن زمان زیاد برای یافتن طول مناسب

## تابعی بر اساس طول قدم

$$\phi(\alpha_i) = f(\mathbf{x}_i + \alpha_i \mathbf{p}_i), \alpha_i > 0$$



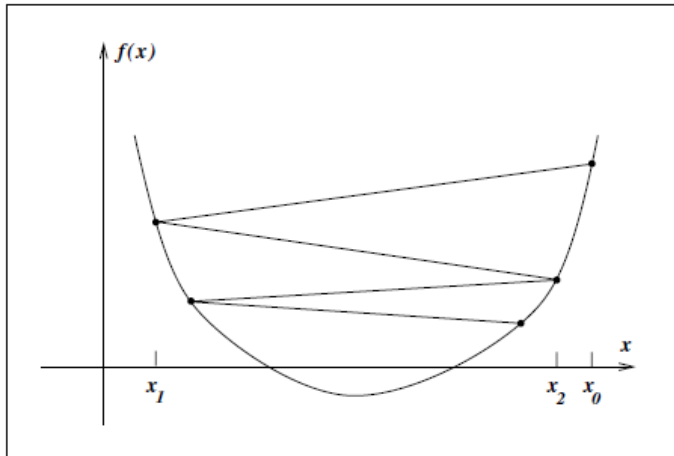
- تابعی تک متغیره
- کمینه‌ساز سراسری
- مشکل یافتن کمینه سراسری
- سیاست عملی‌تر
- جستجو خط‌جزیی

# طول قدم $\alpha$

امر غالب

- الگوریتم جستجو به دنبال یافتن مقداری برای طول قدم
- اتمام الگوریتم هنگام یافتن مقدار برآورده کننده چند شرط
- انجام جستجو خط در دو مرحله
  - کروشه گذاری
  - یافتن بازه مقادیر مطلوب
  - درونیابی/نیمه سازی
  - یافتن مقدار مناسب در بازه مذکور
- ابتدا بررسی شرط خاتمه

# طول قدم $\alpha$



ساده‌ترین شرط

▪ کاهش در مقدار تابع

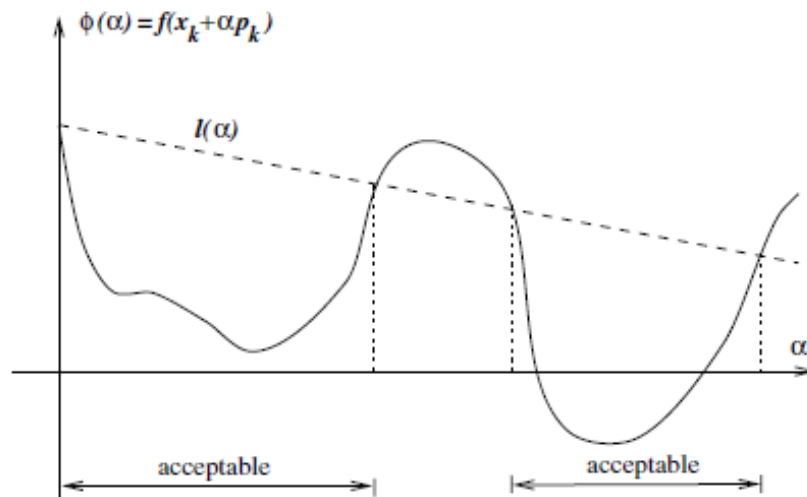
$$f(\mathbf{x}_k + \alpha_i \mathbf{p}_k) < f(\mathbf{x}_k) \quad \square$$

▪ ناکافی جهت هم‌گرایی به کمینه‌ساز

نیاز به شرط کافی کاهش

# شروط وولف

طول قدم در هر مرحله صادق در کاهش کافی تابع هدف



کاهش کافی

$$f(x_i + \alpha_i p_i) \leq f(x_i) + c_1 \alpha_i \nabla f_i^T p_i$$

$$c_1 \in (0, 1)$$

تناسب کاهش  $f$  با

طول قدم  $\alpha_i$

مشتق جهت‌دار  $\nabla f_i^T p_i$

$$l(\alpha_i) = f(x_i) + c_1 \alpha_i \nabla f_i^T p_i$$

شهره به شرط ارمیخو

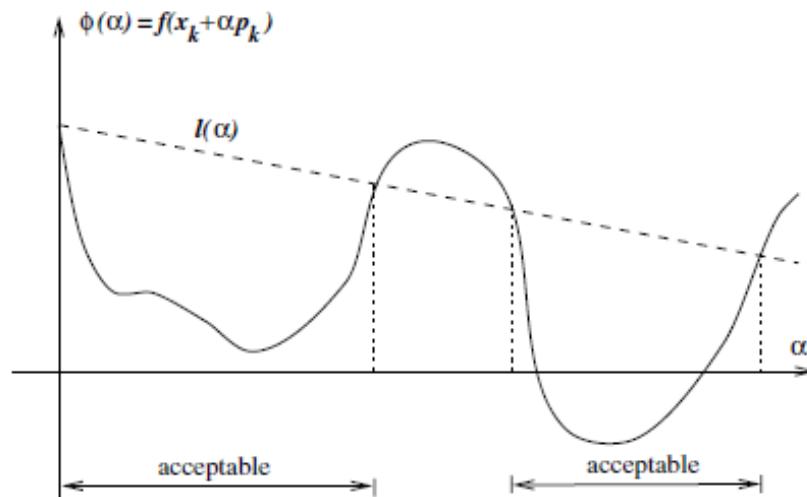
$$\phi(\alpha_i) \leq l(\alpha_i)$$

اختصاص مقداری کوچک به  $c_1$

$$10^{-4}$$

# شروط وولف

طول قدم در هر مرحله صادق در کاهش کافی تابع هدف



کاهش کافی

$$f(x_i + \alpha_i p_i) \leq f(x_i) + c_1 \alpha_i \nabla f_i^T p_i$$

$$c_1 \in (0, 1)$$

تناسب کاهش  $f$  با

طول قدم  $\alpha_i$

مشتق جهت‌دار  $\nabla f_i^T p_i$

$$l(\alpha_i) = f(x_i) + c_1 \alpha_i \nabla f_i^T p_i$$

شهره به شرط ارمینگو

$$\phi(\alpha_i) \leq l(\alpha_i)$$

اختصاص مقداری کوچک به  $c_1$

اما «شرط کافی» مکفی نیست!

# شروط وولف

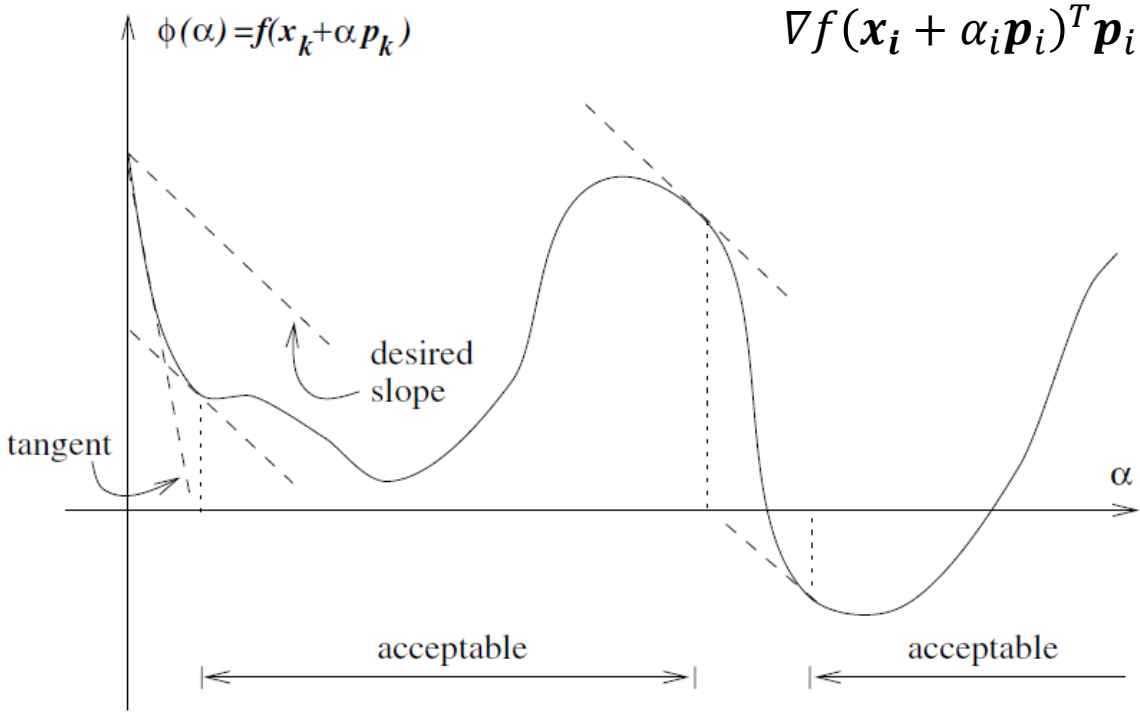
جهت حذف مقادیر اقدام کوچک  
 ▪ شرط انحنا

$$\nabla f(\mathbf{x}_i + \alpha_i \mathbf{p}_i)^T \mathbf{p}_i \geq c_2 \alpha_i \nabla f(\mathbf{x}_i)^T \mathbf{p}_i$$

$$c_2 \in (c_1, 1)$$

$$\phi'(\alpha_i) = \nabla f(\mathbf{x}_i + \alpha_i \mathbf{p}_i)^T \mathbf{p}_i$$

- شیب  $\phi$  در  $\alpha_i$  بزرگتر از مضربی از شیب ابتدایی  $\phi'(0)$  (یا بزرگتر از  $c_2 \phi'(0)$ )
- شیب خیلی منفی
- امکان کاهش بیشتر
- شیب کمی منفی یا کمی مثبت
- ته خط!
- عدم انتظار کاهش بیشتر در مقدار تابع



# شروط وولف

جهت حذف مقادیر اقدام کوچک  
▪ شرط انحنای

$$\nabla f(\mathbf{x}_i + \alpha_i \mathbf{p}_i)^T \mathbf{p}_i \geq c_2 \alpha_i \nabla f(\mathbf{x}_i)^T \mathbf{p}_i$$

$$c_2 \in (c_1, 1) \quad \square$$

$$\phi'(\alpha_i) = \nabla f(\mathbf{x}_i + \alpha_i \mathbf{p}_i)^T \mathbf{p}_i \quad \square$$

▪ شیب  $\phi$  در  $\alpha_i$  بزرگتر از مضربی از شیب ابتدایی  $\phi'(0)$  (یا بزرگتر از  $c_2 \phi'(0)$ )

▪ شیب خیلی منفی

▪ امکان کاهش بیشتر

▪ شیب کمی منفی یا کمی مثبت

▪ ته خط!

▪ عدم انتظار کاهش بیشتر در مقدار تابع

▪  $c_2$

▪ روش نیوتن یا شبه نیوتن معمولاً برابر ۰٫۹

▪ روش گرادیان مزدوج معمولاً برابر ۰٫۱



# شروط وولف

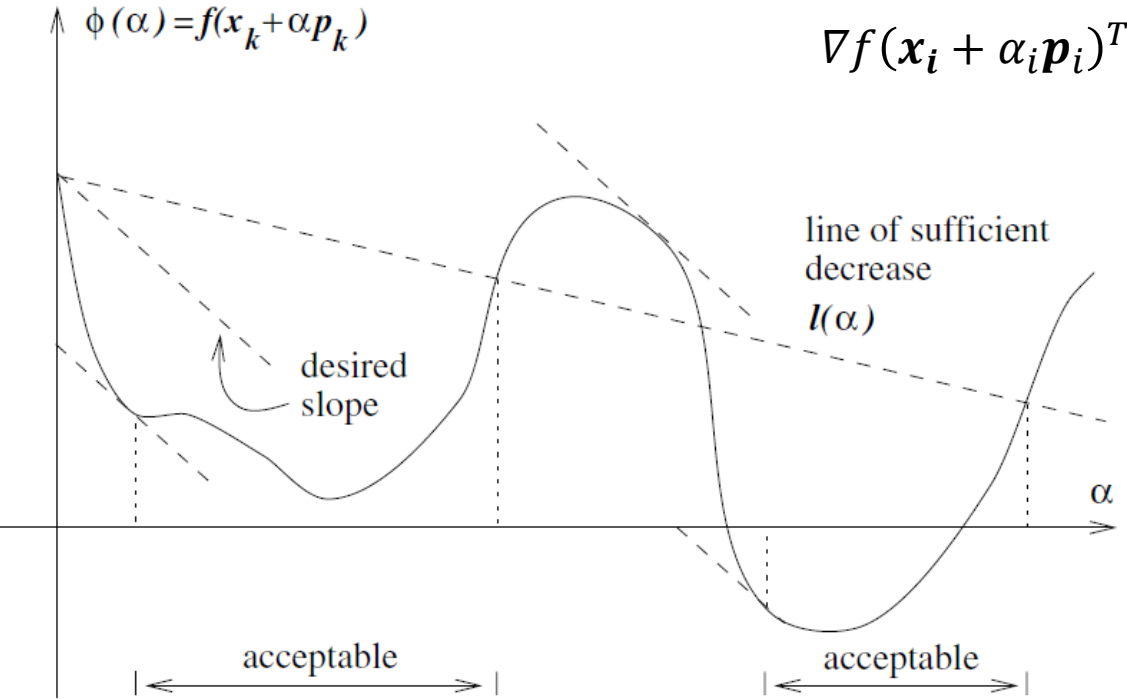
▪ شرط ارميخو

$$f(\mathbf{x}_i + \alpha_i \mathbf{p}_i) \leq f(\mathbf{x}_i) + c_1 \alpha_i \nabla f(\mathbf{x}_i)^T \mathbf{p}_i$$

▪ شرط انحناء

$$\nabla f(\mathbf{x}_i + \alpha_i \mathbf{p}_i)^T \mathbf{p}_i \geq c_2 \alpha_i \nabla f(\mathbf{x}_i)^T \mathbf{p}_i$$

$$0 < c_1 < c_2 < 1$$



# شروط قوی وولف

امکان وجود شرایط وولف بدون اینکه طول مناسبی باشد

شروط قوی وولف

$$f(\mathbf{x}_i + \alpha_i \mathbf{p}_i) \leq f(\mathbf{x}_i) + c_1 \alpha_i \nabla f(\mathbf{x}_i)^T \mathbf{p}_i$$

$$|\nabla f(\mathbf{x}_i + \alpha_i \mathbf{p}_i)^T \mathbf{p}_i| \leq c_2 |\alpha_i \nabla f(\mathbf{x}_i)^T \mathbf{p}_i|$$

$$0 < c_1 < c_2 < 1 \quad \blacksquare$$

# شروط قوی وولف

امکان وجود شرایط وولف بدون اینکه طول مناسبی باشد

شروط قوی وولف

$$f(\mathbf{x}_i + \alpha_i \mathbf{p}_i) \leq f(\mathbf{x}_i) + c_1 \alpha_i \nabla f(\mathbf{x}_i)^T \mathbf{p}_i$$

$$|\nabla f(\mathbf{x}_i + \alpha_i \mathbf{p}_i)^T \mathbf{p}_i| \leq c_2 |\alpha_i \nabla f(\mathbf{x}_i)^T \mathbf{p}_i|$$

$$0 < c_1 < c_2 < 1 \quad \blacksquare$$

▪ جلوگیری از مقادیر خیلی مثبت

# شروط گلدشتین

اطمینان از دستیابی طول به کاهش کافی و جلوگیری از کوتاه قدم برداشتن

$$f(\mathbf{x}_i) + (1 - c)\alpha_i \nabla f(\mathbf{x}_i)^T \mathbf{p}_i \leq f(\mathbf{x}_i + \alpha_i \mathbf{p}_i) \leq f(\mathbf{x}_i) + c\alpha_i \nabla f(\mathbf{x}_i)^T \mathbf{p}_i$$

$$0 < c < 0.5$$

همان شرط کاهش کافی

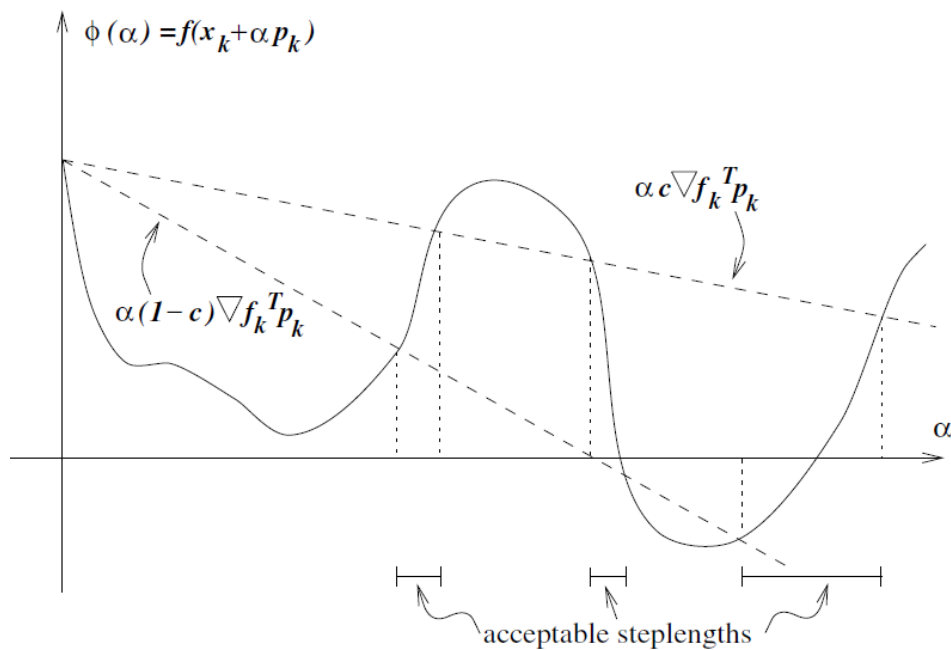
# شروط گلدشتین

اطمینان از دستیابی طول به کاهش کافی و جلوگیری از کوتاه قدم برداشتن

$$f(\mathbf{x}_i) + (1 - c)\alpha_i \nabla f(\mathbf{x}_i)^T \mathbf{p}_i \leq f(\mathbf{x}_i + \alpha_i \mathbf{p}_i) \leq f(\mathbf{x}_i) + c\alpha_i \nabla f(\mathbf{x}_i)^T \mathbf{p}_i$$

$$0 < c < 0.5$$

کنترل طول قدم از پائین



# کاهش کافی با پس روی

صرفاً استفاده از شرط کافی و نه شرط دوم

الگوریتم جستجو خط با پس روی

- انتخاب  $\bar{\alpha} > 0$  و  $\rho \in (0,1)$  و  $c \in (0,1)$  و  $\alpha = \bar{\alpha}$
- تازمان  $f(\mathbf{x}_i + \alpha \mathbf{p}_i) \leq f(\mathbf{x}_i) + c\alpha \nabla f(\mathbf{x}_i)^T \mathbf{p}_i$
- }
- $\alpha = \rho \alpha$
- {
- $\alpha_i = \alpha$

امکان تخصیص پویای  $\rho$  در هر مرحله

# طول قدم $\alpha$

طول قدم

اندازه هر پرش

▪ بزرگ

▪ همگرا نشدن

▪ کوچک

▪ کند

یکی از راه حل ها

▪ شروع با طول قدم بزرگ

▪ مثلاً ۱

▪ کاهش طول قدم هنگام اضافه رفتن  $\alpha = 0.8\alpha$

▪ کاهش دوباره مقدار در صورت کافی نبودن کاهش مقدار قبلی

▪ ادامه تا یافتن طول قدم مناسب

راه دیگر  $\alpha_i = \frac{c}{i}$

# طول قدم $\alpha$ - ادامه

روش دیگر

- وابسته به شکل تابع
- در طول بهینه‌سازی
- پیگیری گرادیان‌های بدست آمده در گذشته
- آداگراد

$$\mathbf{s} = 0$$

تازمان همگرا نشدن

$$\mathbf{g} = -\nabla f(\mathbf{x}_i)$$

$$\mathbf{s} = \mathbf{s} + \|\nabla f(\mathbf{x}_i)\|^2$$

$$\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha \frac{\mathbf{g}}{\sqrt{\mathbf{s} + \epsilon}}$$



# طول قدم $\alpha$ - ادامه

## مشکل آداگراد

- جمع مربعات در مخرج
- بی‌اندازه کوچک شدن طول قدم
- آدالدا
- جلوگیری از کوچک شدن بی‌اندازه مقدار طول قدم
- کاهش طول و استفاده از قبلی‌ها
- با استفاده از پنجره از چند مقدار قبلی

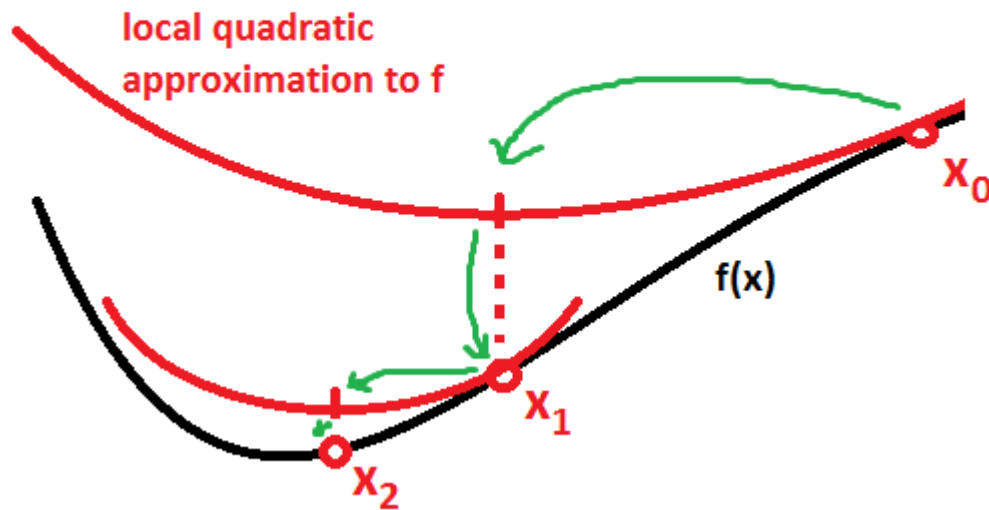
$$\begin{aligned} \mathbf{s} &= 0 \\ \text{تازمان همگرا نشدن} \\ \mathbf{g} &= -\nabla f(\mathbf{x}_i) \\ \mathbf{s} &= \gamma \mathbf{s} + (1 - \gamma) \|\nabla f(\mathbf{x}_i)\|^2 \\ \mathbf{x}_{i+1} &= \mathbf{x}_i + \alpha \frac{\mathbf{g}}{\sqrt{\mathbf{s} + \epsilon}} \end{aligned}$$

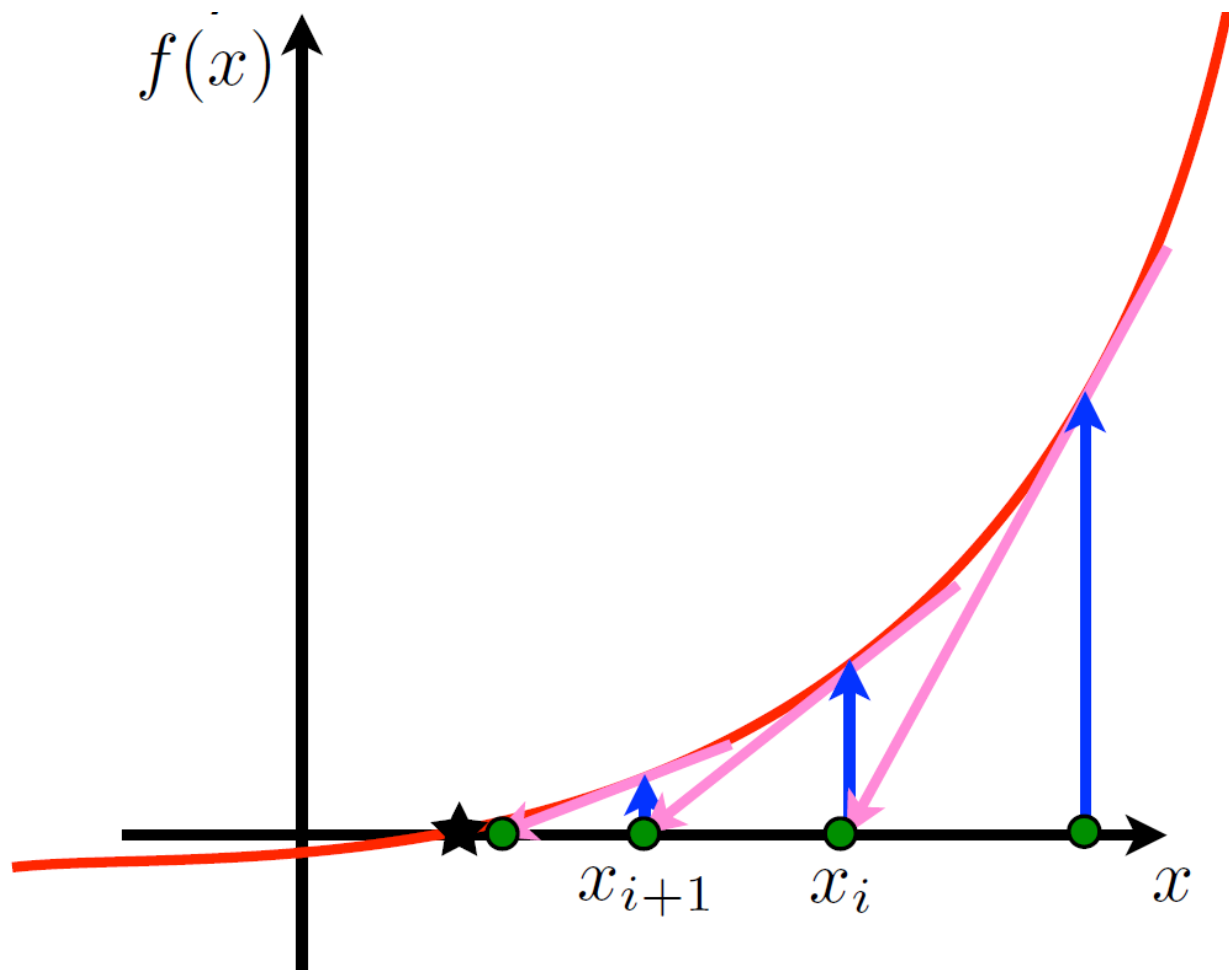
# روش نیوتن

گرادیان نزولی مبتنی بر مشتق مرتبه اول

$$x_{i+1} = x_i - \alpha_i (\nabla f(x_i)) \Leftarrow \text{گن}$$

روشی مبتنی بر مرتبه دوم  
▪ ماتریس هسی





# روش نیوتن

گرادیان نزولی مبتنی بر مشتق مرتبه اول

$$x_{i+1} = x_i - \alpha_i (\nabla f(x_i)) \Leftarrow \text{گن}$$

روشی مبتنی بر مرتبه دوم  
▪ ماتریس هسی

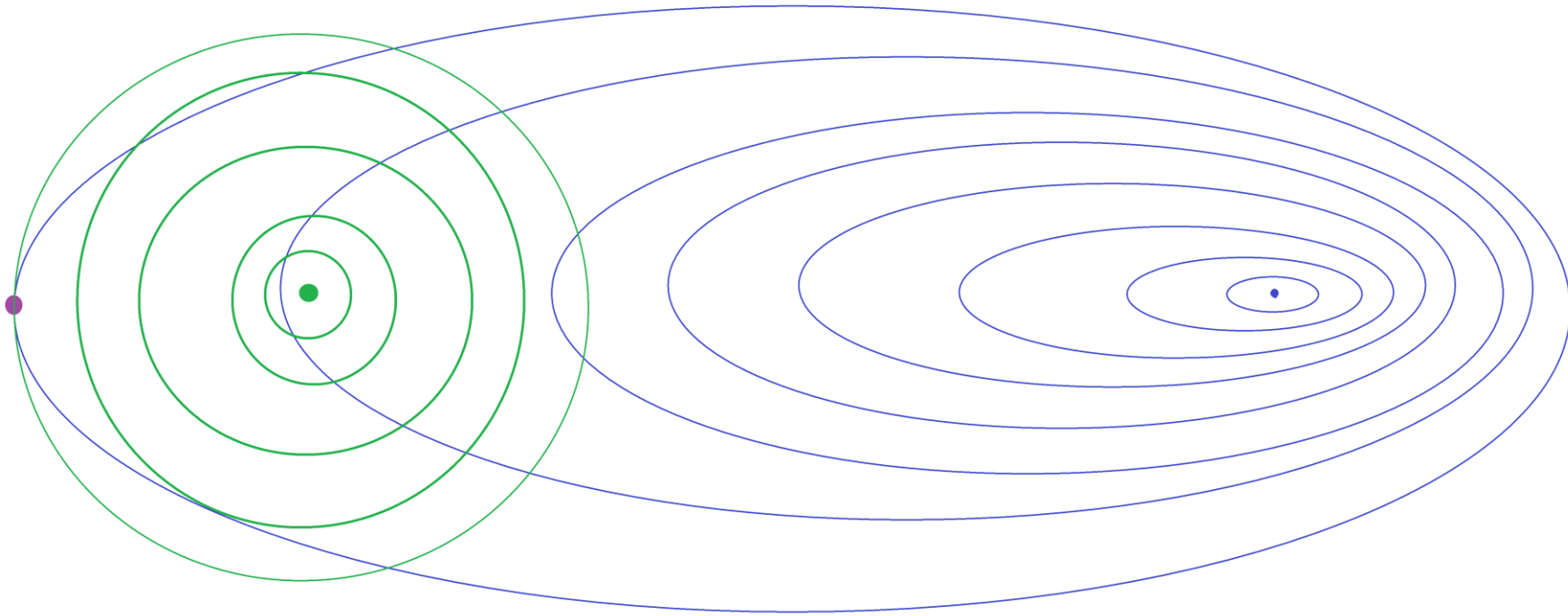
یادآوری جهت تقریب ذهن- برای یافتن صفر (یک بعدی)

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)}$$

کمینه یا بیشینه  
▪ به دنبال اینکه مشتق صفر باشد

$$x_{i+1} = x_i - \frac{f'(x_i)}{f''(x_i)}$$

# روش نیوتن - ادامه



# روش نیوتن - ادامه

$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$

▪ به دنبال یافتن کمینه

قضیه تیلور جهت  $\mathbf{x} \in \mathbb{R}^n$  نزدیک  $\mathbf{a}$

$$f(\mathbf{x}) \approx f(\mathbf{a}) + \nabla f(\mathbf{a})^T (\mathbf{x} - \mathbf{a}) + \frac{1}{2} (\mathbf{x} - \mathbf{a})^T \nabla^2 f(\mathbf{a}) (\mathbf{x} - \mathbf{a})$$

$$\mathbf{g} = \nabla f(\mathbf{a})$$

$$H = \nabla^2 f(\mathbf{a}) = \frac{\partial^2 y}{\partial x_i \partial x_j}, i, j \in \{1, \dots, n\}$$

دیگر روش نوشتن

$$f(\mathbf{x} + \mathbf{p}) \approx f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T \nabla^2 f(\mathbf{x}) \mathbf{p}$$

## روش نیوتن - ادامه

$$f(\mathbf{x}_k + \mathbf{p}) \approx f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T \nabla^2 f(\mathbf{x}_k) \mathbf{p}$$

## روش نیوتن - ادامه

$$f(\mathbf{x}_k + \mathbf{p}) \approx f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T \nabla^2 f(\mathbf{x}_k) \mathbf{p}$$
$$f(\mathbf{x}_k + \mathbf{p}) \approx f(\mathbf{x}_k) + \mathbf{g}_k^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T \mathbf{H}_k \mathbf{p}$$



## روش نیوتن - ادامه

$$f(\mathbf{x}_k + \mathbf{p}) \approx f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T \nabla^2 f(\mathbf{x}_k) \mathbf{p}$$
$$f(\mathbf{x}_k + \mathbf{p}) \approx f(\mathbf{x}_k) + \underbrace{\mathbf{g}_k^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T \mathbf{H}_k \mathbf{p}}_{q(\mathbf{p})}$$

## روش نیوتن - ادامه

$$f(\mathbf{x}_k + \mathbf{p}) \approx f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T \nabla^2 f(\mathbf{x}_k) \mathbf{p}$$
$$f(\mathbf{x}_k + \mathbf{p}) \approx f(\mathbf{x}_k) + \underbrace{\mathbf{g}_k^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T \mathbf{H}_k \mathbf{p}}_{q(\mathbf{p})}$$

روش نیوتن:  $\mathbf{p}$  کاهنده تقریب درجه دو  $q(\mathbf{p})$

## روش نیوتن - ادامه

$$f(\mathbf{x}_k + \mathbf{p}) \approx f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T \nabla^2 f(\mathbf{x}_k) \mathbf{p}$$
$$f(\mathbf{x}_k + \mathbf{p}) \approx f(\mathbf{x}_k) + \underbrace{\mathbf{g}_k^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T \mathbf{H}_k \mathbf{p}}_{q(\mathbf{p})}$$

روش نیوتن:  $\mathbf{p}$  کاهنده تقریب درجه دو  $q(\mathbf{p})$

$$\nabla q(\mathbf{p}) = \mathbf{g} + \mathbf{H}\mathbf{p} = \mathbf{0}$$

## روش نیوتن - ادامه

$$f(\mathbf{x}_k + \mathbf{p}) \approx f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T \nabla^2 f(\mathbf{x}_k) \mathbf{p}$$
$$f(\mathbf{x}_k + \mathbf{p}) \approx f(\mathbf{x}_k) + \underbrace{\mathbf{g}_k^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T \mathbf{H}_k \mathbf{p}}_{q(\mathbf{p})}$$

روش نیوتن:  $\mathbf{p}$  کاهنده تقریب درجه دو  $q(\mathbf{p})$

$$\nabla q(\mathbf{p}) = \mathbf{g} + \mathbf{H}\mathbf{p} = \mathbf{0}$$

$$\mathbf{H}\mathbf{p} = -\mathbf{g}$$

## روش نیوتن - ادامه

$$f(\mathbf{x}_k + \mathbf{p}) \approx f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T \nabla^2 f(\mathbf{x}_k) \mathbf{p}$$
$$f(\mathbf{x}_k + \mathbf{p}) \approx f(\mathbf{x}_k) + \underbrace{\mathbf{g}_k^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T \mathbf{H}_k \mathbf{p}}_{q(\mathbf{p})}$$

روش نیوتن:  $\mathbf{p}$  کاهنده تقریب درجه دو  $q(\mathbf{p})$

$$\nabla q(\mathbf{p}) = \mathbf{g} + \mathbf{H}\mathbf{p} = \mathbf{0}$$

معادله نیوتن

$$\mathbf{H}\mathbf{p} = -\mathbf{g}$$

# روش نیوتن - ادامه

$$f(\mathbf{x}_k + \mathbf{p}) \approx f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T \nabla^2 f(\mathbf{x}_k) \mathbf{p}$$
$$f(\mathbf{x}_k + \mathbf{p}) \approx f(\mathbf{x}_k) + \underbrace{\mathbf{g}_k^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T \mathbf{H}_k \mathbf{p}}_{q(\mathbf{p})}$$

روش نیوتن:  $\mathbf{p}$  کاهنده تقریب درجه دو  $q(\mathbf{p})$

$$\nabla q(\mathbf{p}) = \mathbf{g} + \mathbf{H}\mathbf{p} = \mathbf{0}$$

معادله نیوتن

$$\mathbf{H}\mathbf{p} = -\mathbf{g}$$
$$\mathbf{p}_k = -\mathbf{H}_k^{-1} \mathbf{g}_k$$

$\mathbf{p}_k$  جهت نیوتن

# روش نیوتن - ادامه

$$\mathbf{p}_k = -\mathbf{H}_k^{-1} \mathbf{g}_k$$

با جستجو خط  $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k$

# روش نیوتن - ادامه

الگوریتم روش نیوتن

مقداردهی اولیه  $\mathbf{x}_0 \in \mathbb{R}^n$

تکرار  $\mathbf{x}_{i+1} = \mathbf{x}_i - H^{-1} g$

$$g = \nabla f(\mathbf{x}_i) \cdot$$

$$H = \nabla^2 f(\mathbf{x}_i) \cdot$$



# روش نیوتن - حل معادله غیرخطی

$$\min_{x \in \mathbb{R}^n} f(x)$$
$$\nabla f(x) = \mathbf{g}(x) = \mathbf{0}$$

نیاز به حل دستگاه n-معادله غیرخطی  $\mathbf{g}_i(x) = 0$

$$g(x_k + p) \approx g(x_k) + H(x_k)p$$

# روش نیوتن - حل معادله غیرخطی

$$\min_{x \in \mathbb{R}^n} f(x)$$
$$\nabla f(x) = \mathbf{g}(x) = \mathbf{0}$$

نیاز به حل دستگاه n-معادله غیرخطی  $\mathbf{g}_i(x) = 0$

$$\mathbf{g}(\mathbf{x}_k + \mathbf{p}) \approx \mathbf{g}(\mathbf{x}_k) + H(\mathbf{x}_k)\mathbf{p} = \mathbf{0}$$

$$\mathbf{p} = H^{-1}(\mathbf{x}_k)\mathbf{g}(\mathbf{x}_k)$$

## روش نیوتن - ادامه

$$\mathbf{p}_k = H^{-1}(\mathbf{x}_k)g(\mathbf{x}_k)$$
$$f'_p(\mathbf{x}_k) = g_k^T \mathbf{p}_k = -g_k^T H_k^{-1} g_k < 0, g_k \neq 0$$

برآورده‌پذیر در صورت مثبت معین بودن  $H$

در صورت مثبت معین نبودن  $H$ : افزودن قطر مثبت  $\Delta_k$  به طوری که

$$H_k + \Delta_k > \epsilon I \Leftrightarrow \lambda_i(H_k) > \epsilon$$

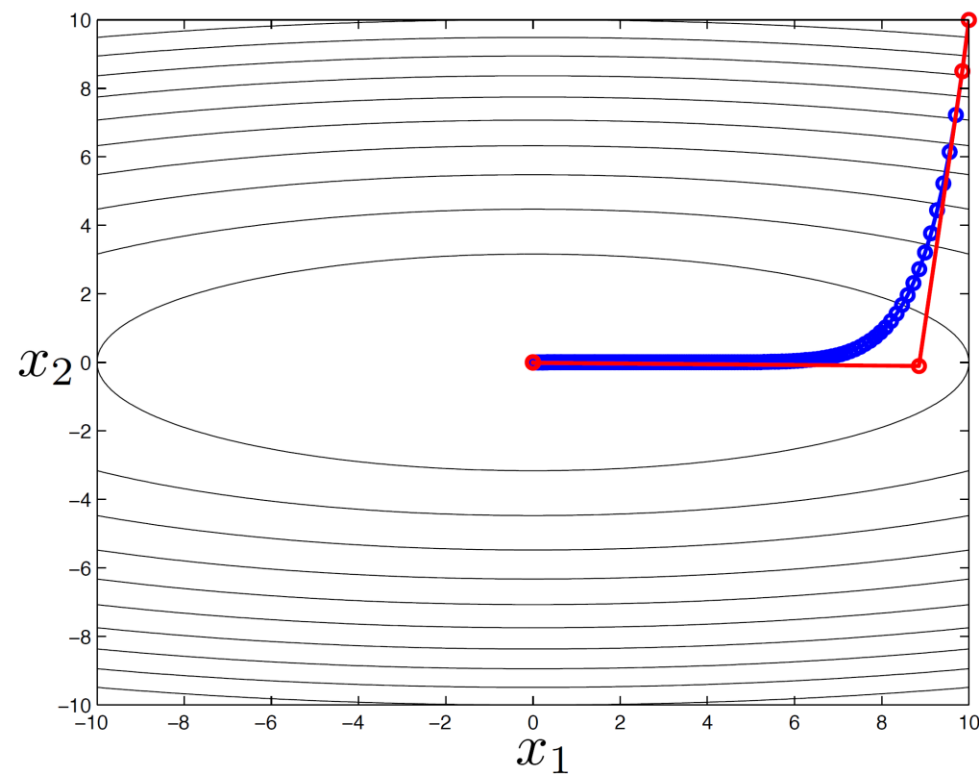
روش نیوتن تغییر یافته

$$\mathbf{p}_k = (H_k + \Delta_k)^{-1} g_k$$

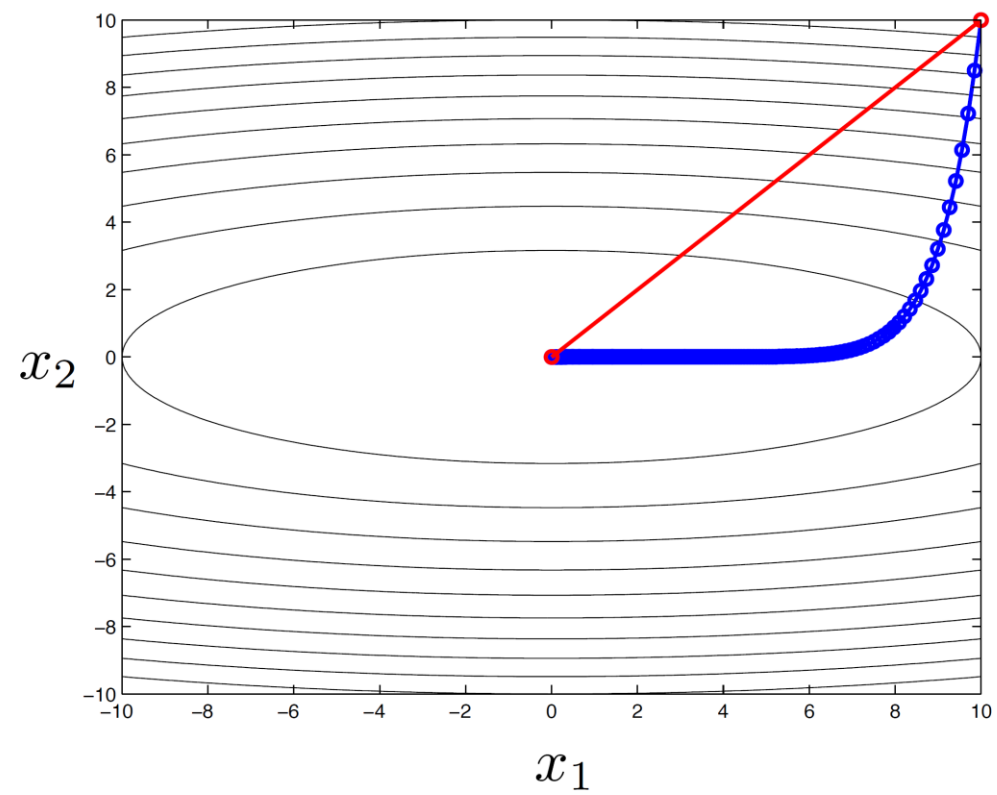
# مثال - تندترین نزول در برابر گرادیان مزدوج

$$f(\mathbf{x}) = x_1^2 + 10x_2^2 \quad \mathbf{A} = \begin{pmatrix} 1 & 0 \\ 0 & 10 \end{pmatrix}, \mathbf{b} = 0$$

$$\alpha_i = 0.015$$



# مثال



# روش نیوتن - ادامه

مزایا

- بسیار سریع

معایب

- ماتریس هسی ممکن است مثبت معین نباشد

- تغییر به گرادیان نزولی

- به جای یافتن معکوس ماتریس هسی، حل  $\mathbf{H}\mathbf{y}=\mathbf{g}$  برای  $\mathbf{y}$

- استفاده از  $\mathbf{x}_{i+1} = \mathbf{x}_i - \mathbf{y}$

- $\mathbf{x}_{i+1} = \mathbf{x}_i - \alpha\mathbf{y}$

- بیشتر اوقات کار نمی کند

- کار می کند وقتی که خیلی نزدیک کمینه باشد

- پیشنهاد ترکیب گرادیان نزولی و سپس روش نیوتن

# منابع

[نازهدل]

What are the differences between the different gradient-based numerical optimization methods?, <https://scicomp.stackexchange.com/questions/26960/what-are-the-differences-between-the-different-gradient-based-numerical-optimiza>

“An overview of gradient descent optimization algorithms” <https://ruder.io/optimizing-gradient-descent/>, 2016