

Lecture 16: Gradient descent – modifications

Ciaran Evans

Recall: gradient descent

- ▶ **Goal:** Minimize $f(\mathbf{x})$
- ▶ Iterative updates: $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \nabla f(\mathbf{x}^{(k)})$
- ▶ Choosing step size α_k :
 - ▶ one option is backtracking line search with sufficient decrease condition

Limitations of gradient descent

Motivating example: Data on med school admissions for 55 students

- ▶ GPA: student's undergraduate GPA
- ▶ MCAT: student's MCAT score

Function to minimize:

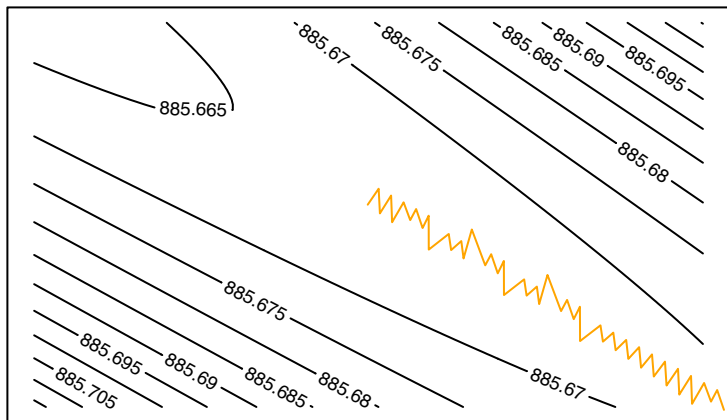
$$f(\beta_0, \beta_1) = \sum_{i=1}^n (\text{MCAT}_i - \beta_0 - \beta_1 \text{GPA}_i)^2$$

- ▶ Gradient descent with backtracking linear search beginning at (0, 0): 6517 iterations

Question: from the activity last time, why does gradient descent need so many iterations here?

Limitations of gradient descent

Gradient descent struggles to traverse long, narrow valleys



Question: What do you notice about the gradient descent path?

Zig-zags!

Limitations of gradient descent

We should expect gradient descent to have a zig-zag path; this makes long, narrow valleys slow.

- ▶ $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \nabla f(\mathbf{x}^{(k)})$
- ▶ Suppose α_k is chosen to minimize

$$f(\mathbf{x}^{(k)} - \alpha_k \nabla f(\mathbf{x}^{(k)}))$$



Claim: $\nabla f(\mathbf{x}^{(k+1)})$ is **orthogonal** (perpendicular) to $\nabla f(\mathbf{x}^{(k)})$

i.e. $\mathbf{x}^{(k)} - \alpha_k \nabla f(\mathbf{x}^{(k)})$

minimizes f along
direction $-\nabla f(\mathbf{x}^{(k)})$

$$\Rightarrow \mathbf{D}_{-\nabla f(\mathbf{x}^{(k)})} f(\mathbf{x}^{(k)}) = 0$$

$$\Rightarrow \nabla f(\mathbf{x}^{(k+1)})^T \nabla f(\mathbf{x}^{(k)}) = 0$$

$$\Rightarrow \nabla f(\mathbf{x}^{(k+1)}) \perp \nabla f(\mathbf{x}^{(k)})$$

increase f
increase f
 $\mathbf{x}^{(k)} - \alpha_k \nabla f(\mathbf{x}^{(k)})$
minimize f
in direction of $-\nabla f$



Overview: modifications of gradient descent

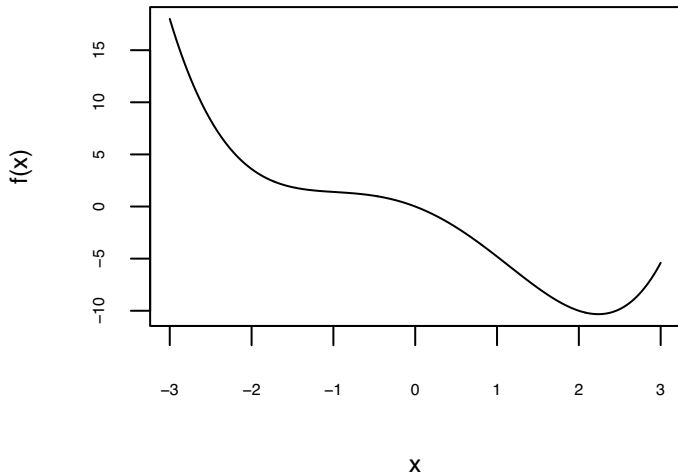
Basic idea: don't always move in the direction of steepest descent

- ▶ Momentum methods: direction of movement is combination of current gradient and previous gradients
- ▶ Subgradient methods: step size is different for each entry in gradient vector
- ▶ Second-order methods: use information about curvature of function (second derivative), not just gradient

It would be impossible to cover all of the possible modifications. My goal is to give you an introduction to some of the common ideas, and their motivation.

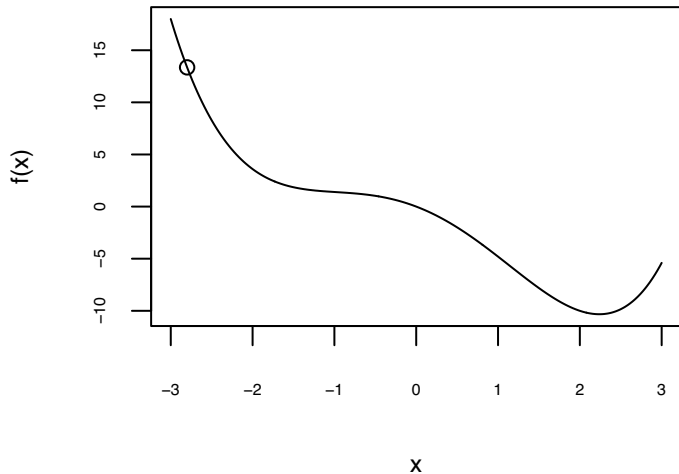
Momentum: motivation

Suppose we wish to minimize this function:



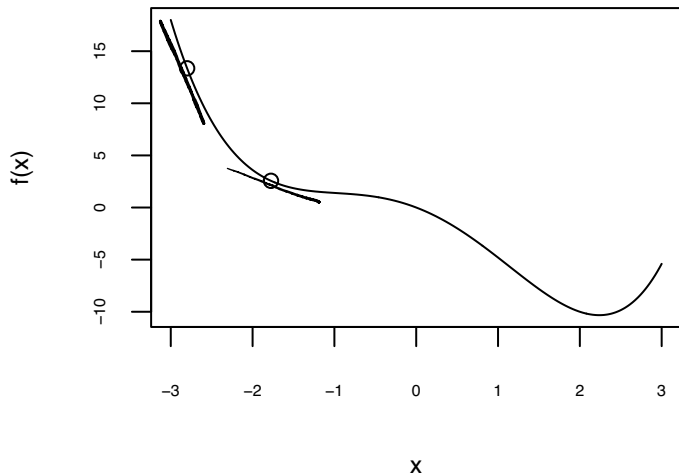
Momentum: motivation

We start at $x^{(0)} = -2.8$:



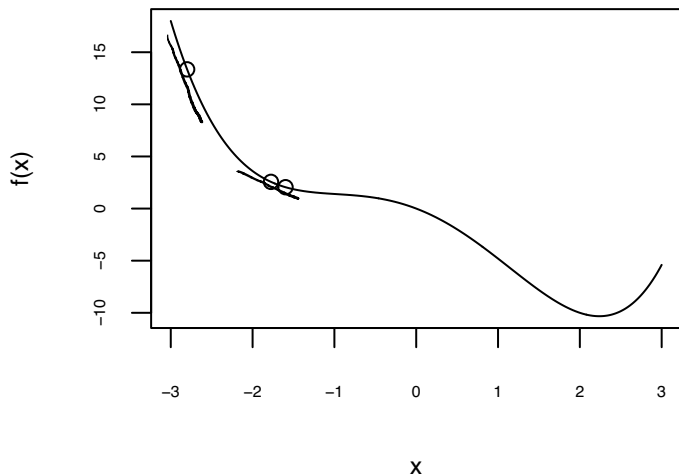
Momentum: motivation

And now we perform gradient descent. After the first iteration:



Momentum: motivation

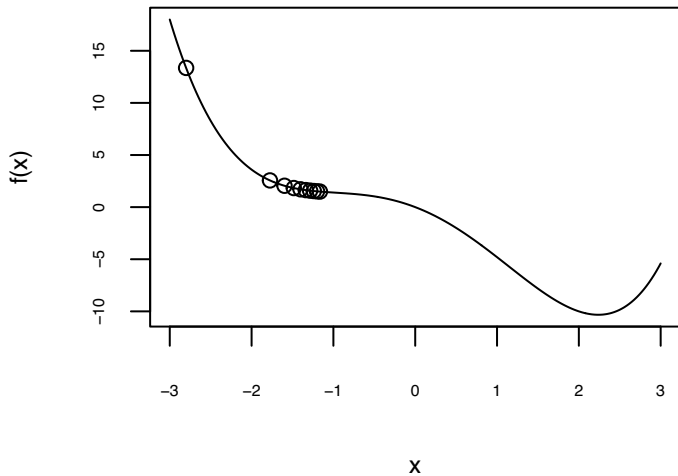
And now we perform gradient descent. After the second iteration:



Question: Why is the second step smaller than the first step?

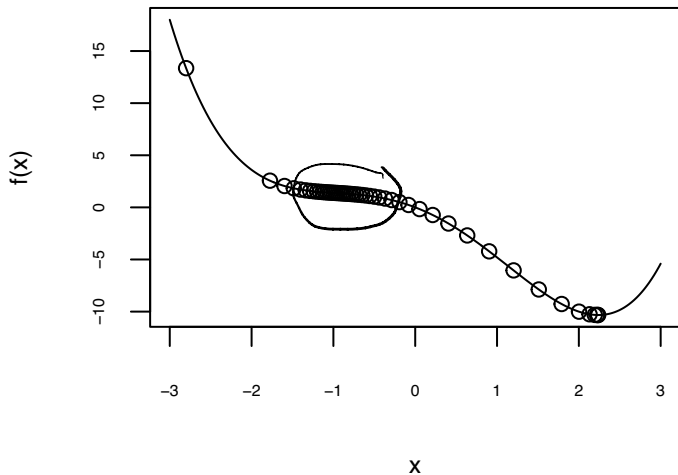
Momentum: motivation

Several more iterations:



Momentum: motivation

Finally, after 50 iterations:



Question: At which part of this function is gradient descent slowest?

Momentum

Gradient descent can be slow in flat areas. Idea: use the momentum from previous gradients.

► Let $0 \leq \beta < 1$ (weight on previous gradients)

► Update rule:

► $\mathbf{x}^{(1)} = \mathbf{x}^{(0)} - \alpha \nabla f(\mathbf{x}^{(0)})$

► $\mathbf{x}^{(2)} = \mathbf{x}^{(1)} - \alpha \nabla f(\mathbf{x}^{(1)}) - \underbrace{\beta \alpha \nabla f(\mathbf{x}^{(0)})}_{\text{previous direction of movement}}$
some weight on previous movement

► $\mathbf{x}^{(3)} = \mathbf{x}^{(2)} - \alpha \nabla f(\mathbf{x}^{(2)}) - \beta \alpha \nabla f(\mathbf{x}^{(1)}) - \beta^2 \alpha \nabla f(\mathbf{x}^{(0)})$
↑
decaying weight on previous gradients

Momentum

Gradient descent can be slow in flat areas. Idea: use the momentum from previous gradients.

- ▶ Let $0 \leq \beta < 1$

- ▶ Update rule:

$$\begin{aligned} \text{▶ } \mathbf{x}^{(1)} &= \mathbf{x}^{(0)} - \underbrace{\alpha \nabla f(\mathbf{x}^{(0)})}_{\mathbf{v}^{(0)}} = \mathbf{x}^{(0)} + \underbrace{\beta \mathbf{v}^{(0)}}_{\text{momentum at first step}} + \mathbf{v}^{(0)} \\ \text{▶ } \mathbf{x}^{(2)} &= \mathbf{x}^{(1)} - \underbrace{\alpha \nabla f(\mathbf{x}^{(1)})}_{\mathbf{v}^{(1)}} - \underbrace{\beta \alpha \nabla f(\mathbf{x}^{(0)})}_{\text{momentum at second step}} = \mathbf{x}^{(1)} + \mathbf{v}^{(1)} \\ \text{▶ } \mathbf{x}^{(3)} &= \mathbf{x}^{(2)} - \underbrace{\alpha \nabla f(\mathbf{x}^{(2)})}_{\mathbf{v}^{(2)}} - \underbrace{\beta \alpha \nabla f(\mathbf{x}^{(1)})}_{\beta \mathbf{v}^{(1)}} - \underbrace{\beta^2 \alpha \nabla f(\mathbf{x}^{(0)})}_{\beta^2 \mathbf{v}^{(0)}} \end{aligned}$$

Written another way:

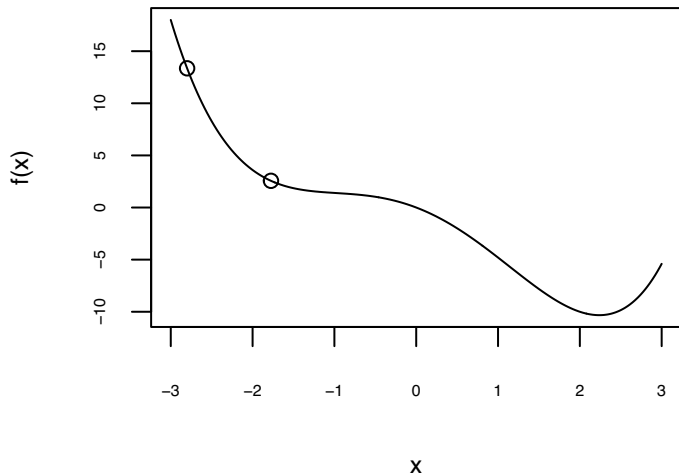
$$\begin{aligned} \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} + \mathbf{v}^{(k)} \quad \leftarrow \text{current momentum} \\ \mathbf{v}^{(k)} &= -\alpha \nabla f(\mathbf{x}^{(k)}) + \beta \mathbf{v}^{(k-1)} \quad \leftarrow \text{previous momentum} \end{aligned}$$

Momentum

- ▶ Let $0 \leq \beta < 1$
- ▶ $\mathbf{v}^{(0)} = -\alpha \nabla f(\mathbf{x}^{(0)})$
- ▶ $\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \mathbf{v}^{(0)}$
- ▶ $\mathbf{x}^{(k)} = -\alpha \nabla f(\mathbf{x}^{(k)}) + \beta \mathbf{v}^{(k-1)}$
- ▶ $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mathbf{v}^{(k)}$

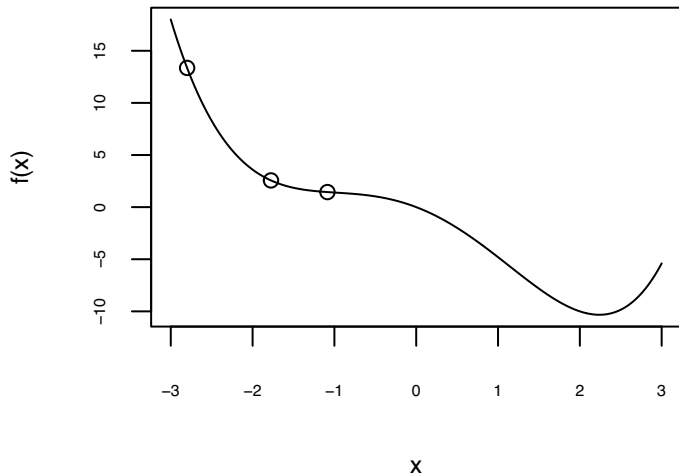
Momentum in action

Starting again at $x^{(0)} = -2.8$. After the first iteration:



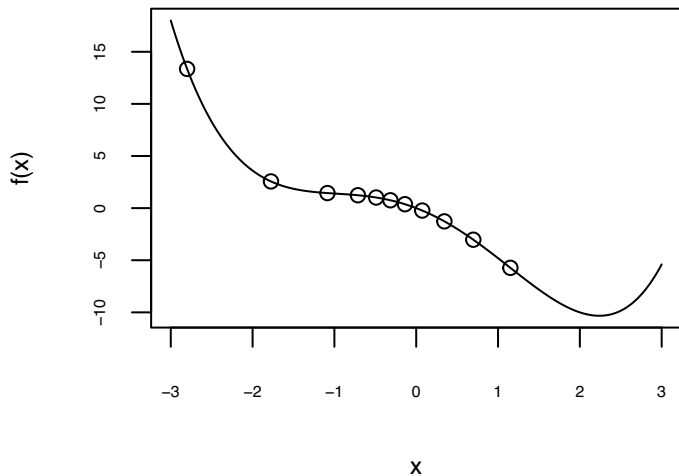
Momentum in action

After the two iterations:



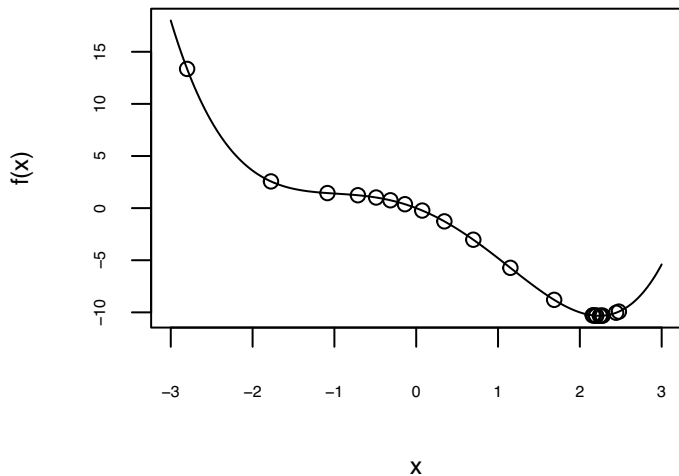
Momentum in action

After 10 iterations:



Momentum in action

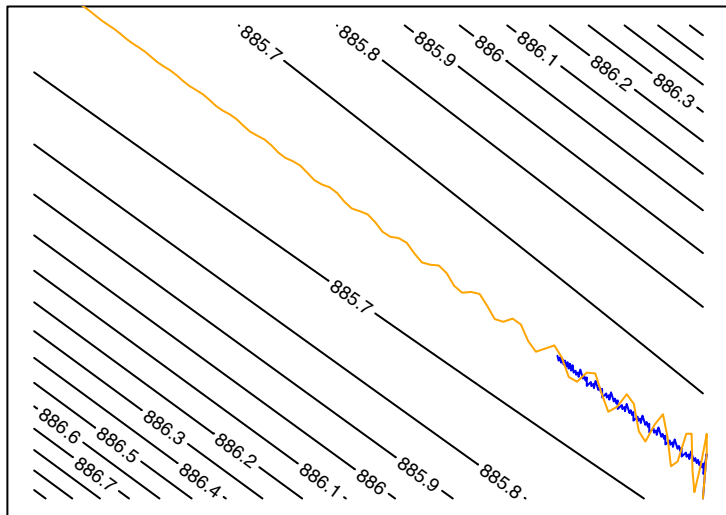
After 20 iterations:



Momentum in action

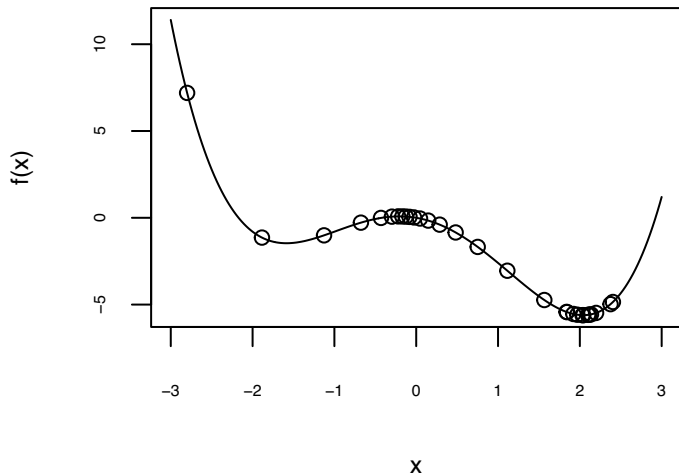
Blue: 100 iterations of gradient descent

Orange: 100 iterations of descent with momentum



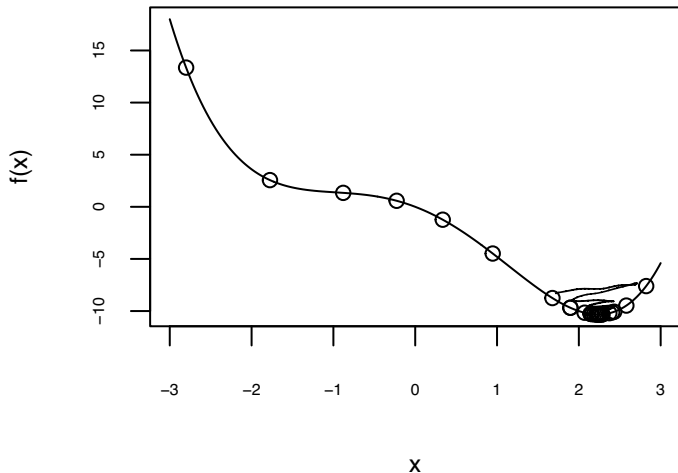
Another example

Momentum can also help overcome local minima:



Drawback of momentum

Question: What do you notice about the iterations near the minimum?



oscillates around

Nesterov momentum

“Heavy ball” momentum:

- ▶ $\mathbf{v}^{(k)} = -\alpha \nabla f(\mathbf{x}^{(k)}) + \beta \mathbf{v}^{(k-1)}$
- ▶ $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mathbf{v}^{(k)}$

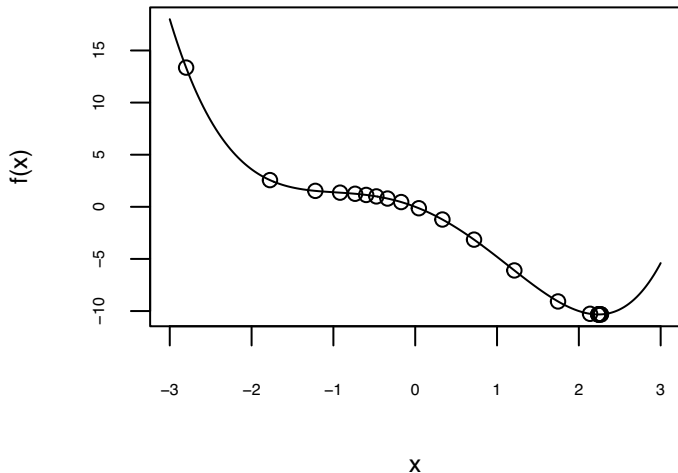
look at gradient in our future
direction of movement, not
just current point

Nesterov momentum:

- ▶ $\mathbf{v}^{(k)} = -\alpha \nabla f(\mathbf{x}^{(k)} + \beta \mathbf{v}^{(k-1)}) + \beta \mathbf{v}^{(k-1)}$
- ▶ $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mathbf{v}^{(k)}$

Nesterov momentum in action

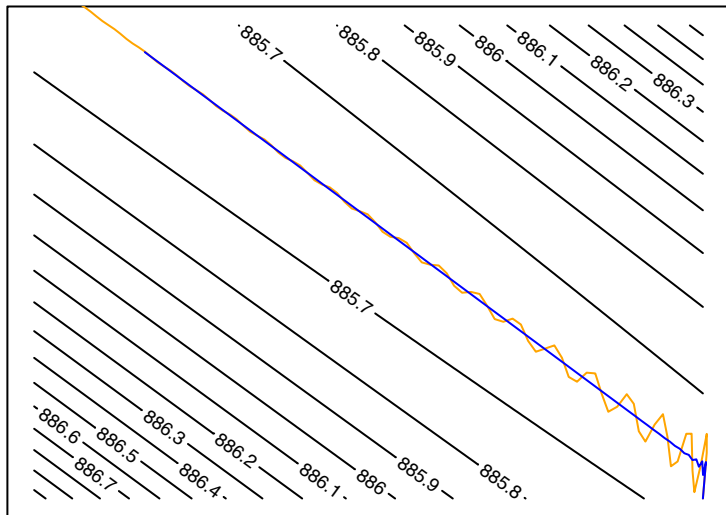
Descent with Nesterov momentum can slow itself down when it gets to the bottom:



Nesterov momentum in action

Blue: Nesterov momentum

Orange: Heavy ball momentum



Subgradient methods

$$\hat{\beta}_0 = 3.9, \hat{\beta}_1 = 9.1$$

Gradient descent: $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \nabla f(\mathbf{x}^{(k)})$

- ▶ The same step size α_k is applied to each element of the gradient $\nabla f(\mathbf{x}^{(k)})$
- ▶ Example gradient from the MedGPA example:

$$\nabla f((4.1, 9.3)) = \begin{pmatrix} 146.1 \\ 521.7 \end{pmatrix} \leftarrow \begin{array}{l} \text{changes in } \beta_1 \\ \text{have a bigger} \\ \text{impact on } f \end{array}$$

Question: Why might we not want to use the same value of α_k for each element of the gradient?

Adagrad

Adagrad (adaptive subgradient) uses a different step size for each element of the gradient:

step size for i th coordinate

$$x_i^{(k+1)} = x_i^{(k)} - \frac{\alpha}{\epsilon + \sqrt{s_i^{(k)}}} \nabla f(x^{(k)})_i$$

i th entry of $x^{(k+1)}$

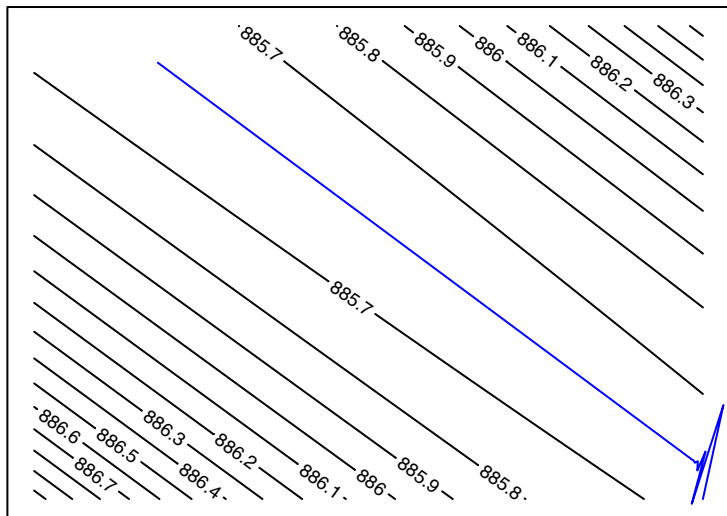
ϵ prevents division by 0

$$s_i^{(k)} = \sum_{j=1}^k (\nabla f(x^{(j)})_i)^2$$

i th entry of $\nabla f(x^{(k)})$

idea: take a smaller step if i th coordinate has a consistently large gradient

Adagrad



Your turn

Practice questions on the course website:

https://sta379-s25.github.io/practice_questions/pq_16.html

- ▶ Try heavy ball momentum
- ▶ Start in class. You are welcome to work with others
- ▶ Practice questions are to help you practice. They are not submitted and not graded
- ▶ Solutions are posted on the course website