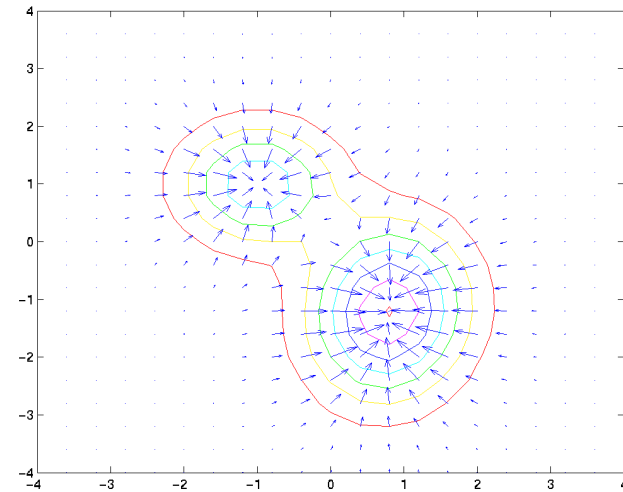
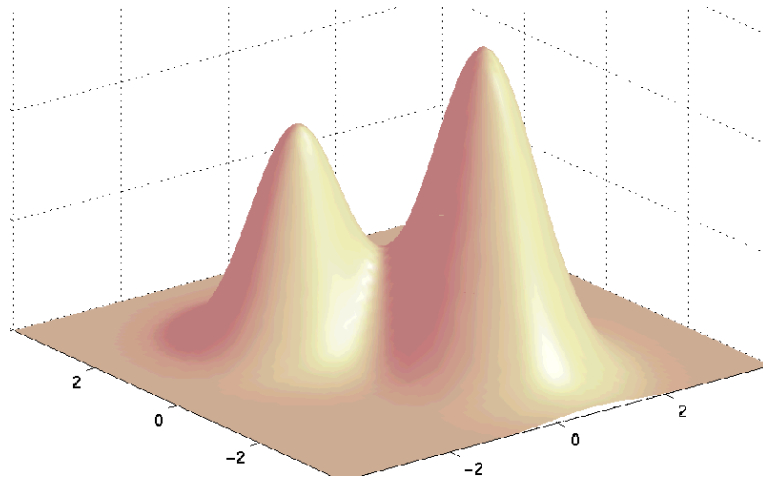


# Advanced Machine Learning

## *Optimisation*

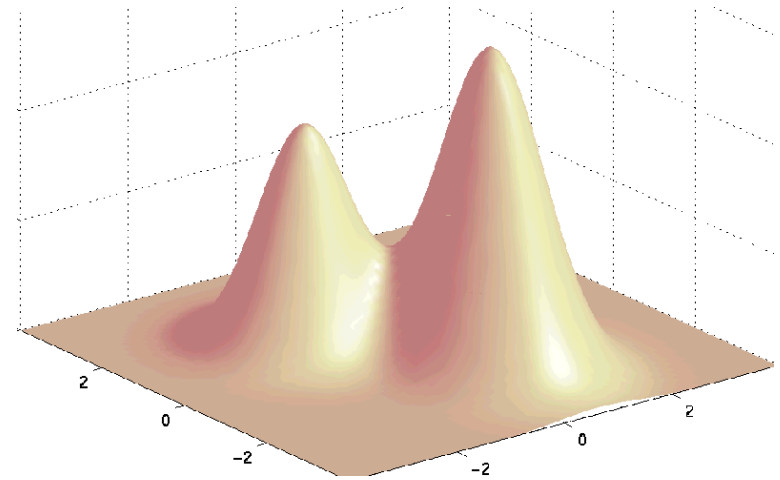


$$z = e^{-(x+1)^2 - (y-1)^2} + 0.6 e^{-(x-1)^2 - 0.5(y+1)^2 + 0.1(x-3)(y-3)}$$

*Gradient descent, quadratic minima, differing length scales*

# Outline

1. **Motivation**
2. Gradient Descent
3. Why Gradient Descent is Difficult



# ML = Optimisation

- Many learning machines can be thought of as functions of the form

$$\hat{y} = f(x|w)$$

(or more generally  $\hat{y} = f(x|w)$ )

- Given an input pattern (set of features)  $x$  the learning machine makes a prediction  $\hat{y}$
- We try to choose the parameters  $w$  so that the predictions are good
- In practice training a learning machine comes down to optimising some loss function

# ML = Optimisation

- Many learning machines can be thought of as functions of the form

$$\hat{y} = f(x|w)$$

(or more generally  $\hat{y} = f(x|w)$ )

- Given an input pattern (set of features)  $x$  the learning machine makes a prediction  $\hat{y}$
- We try to choose the parameters  $w$  so that the predictions are good
- In practice training a learning machine comes down to optimising some loss function

# ML = Optimisation

- Many learning machines can be thought of as functions of the form

$$\hat{y} = f(x|w)$$

(or more generally  $\hat{y} = f(x|w)$ )

- Given an input pattern (set of features)  $x$  the learning machine makes a prediction  $\hat{y}$
- We try to choose the parameters  $w$  so that the predictions are good
- In practice training a learning machine comes down to optimising some loss function

# ML = Optimisation

- Many learning machines can be thought of as functions of the form

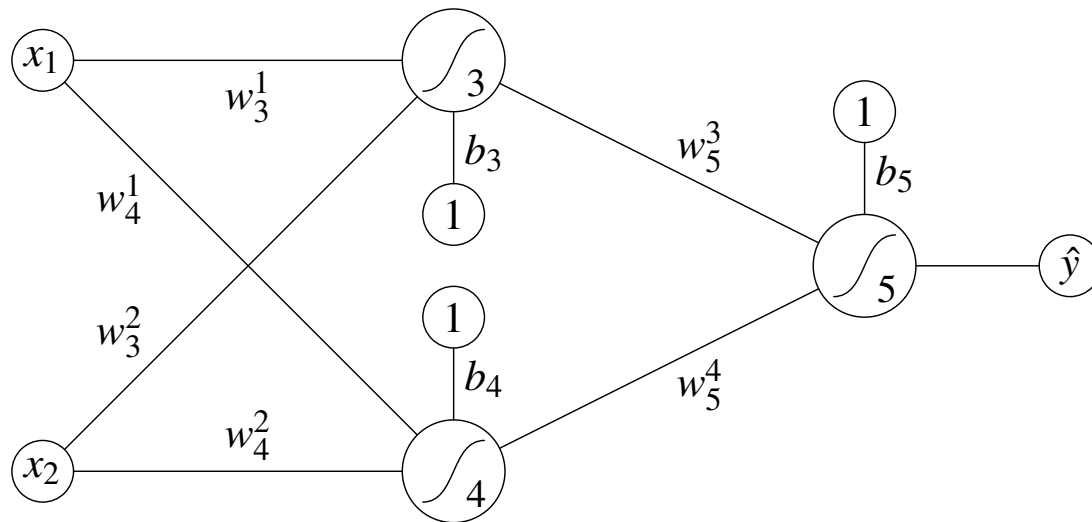
$$\hat{y} = f(x|w)$$

(or more generally  $\hat{y} = f(x|w)$ )

- Given an input pattern (set of features)  $x$  the learning machine makes a prediction  $\hat{y}$
- We try to choose the parameters  $w$  so that the predictions are good
- In practice training a learning machine comes down to optimising some loss function

# MLP

- We can depict a neural network such as an MLP by a diagram



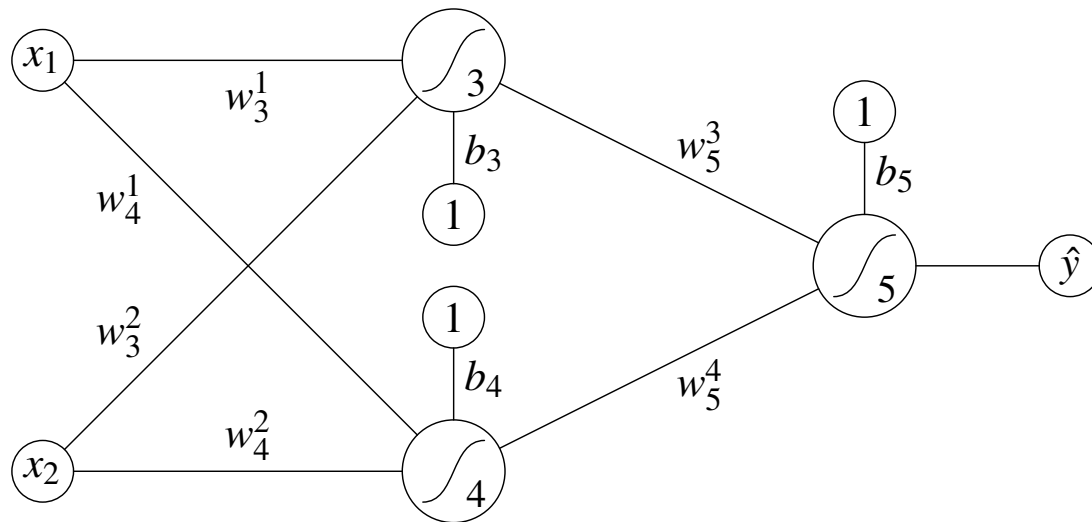
- Stands for the function  $(\hat{y} = f(\mathbf{x}|\mathbf{w}))$

$$\hat{y} = g(w_5^3 g(w_3^1 x_1 + w_3^2 x_2 + b_3) + w_5^4 g(w_4^1 x_1 + w_4^2 x_2 + b_4) + b_5)$$

where, for example,  $g(V) = \frac{1}{1+e^{-V}}$

# MLP

- We can depict a neural network such as an MLP by a diagram



- Stands for the function ( $\hat{y} = f(x|w)$ )

$$\hat{y} = g(w_5^3 g(w_3^1 x_1 + w_3^2 x_2 + b_3) + w_5^4 g(w_4^1 x_1 + w_4^2 x_2 + b_4) + b_5)$$

where, for example,  $g(V) = \frac{1}{1+e^{-V}}$



# Training

- Given a (labelled) training dataset

$$\mathcal{D} = \{(\mathbf{x}_k, y_k) \mid k = 1, \dots, m\}$$

- We define an error or loss function that we want to minimise

$$L(\mathbf{w}|\mathcal{D}) = \frac{1}{m} \sum_{k=1}^m (f(\mathbf{x}_k|\mathbf{w}) - y_k)^2$$

- We then use the machine with the weights  $\mathbf{w}^*$  which minimise  $L(\mathbf{w}|\mathcal{D})$

# Training

- Given a (labelled) training dataset

$$\mathcal{D} = \{(\mathbf{x}_k, y_k) \mid k = 1, \dots, m\}$$

- We define an error or loss function that we want to minimise

$$L(\mathbf{w}|\mathcal{D}) = \frac{1}{m} \sum_{k=1}^m (f(\mathbf{x}_k|\mathbf{w}) - y_k)^2$$

- We then use the machine with the weights  $\mathbf{w}^*$  which minimise  $L(\mathbf{w}|\mathcal{D})$

# Training

- Given a (labelled) training dataset

$$\mathcal{D} = \{(\mathbf{x}_k, y_k) \mid k = 1, \dots, m\}$$

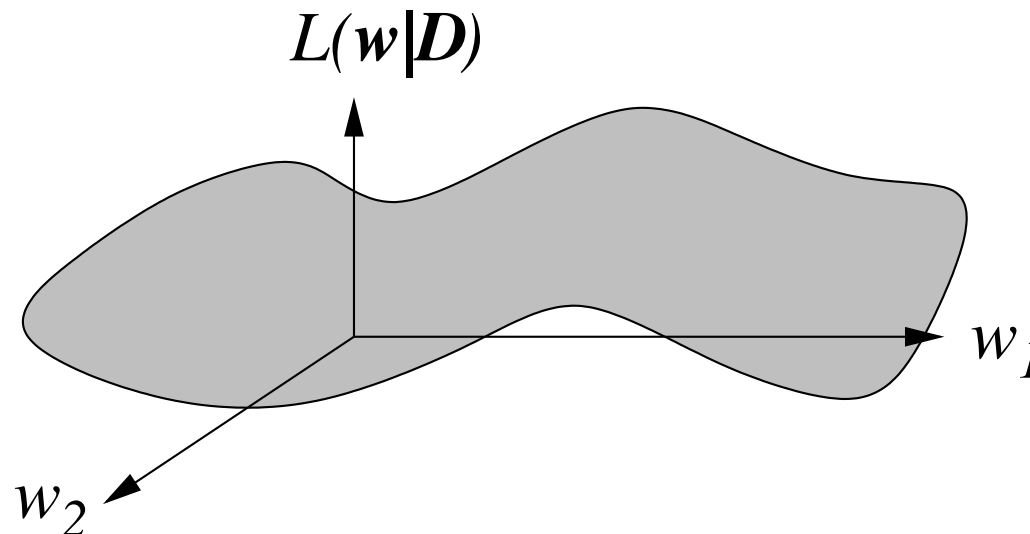
- We define an error or loss function that we want to minimise

$$L(\mathbf{w}|\mathcal{D}) = \frac{1}{m} \sum_{k=1}^m (f(\mathbf{x}_k|\mathbf{w}) - y_k)^2$$

- We then use the machine with the weights  $\mathbf{w}^*$  which minimise  $L(\mathbf{w}|\mathcal{D})$

# Computing Gradients

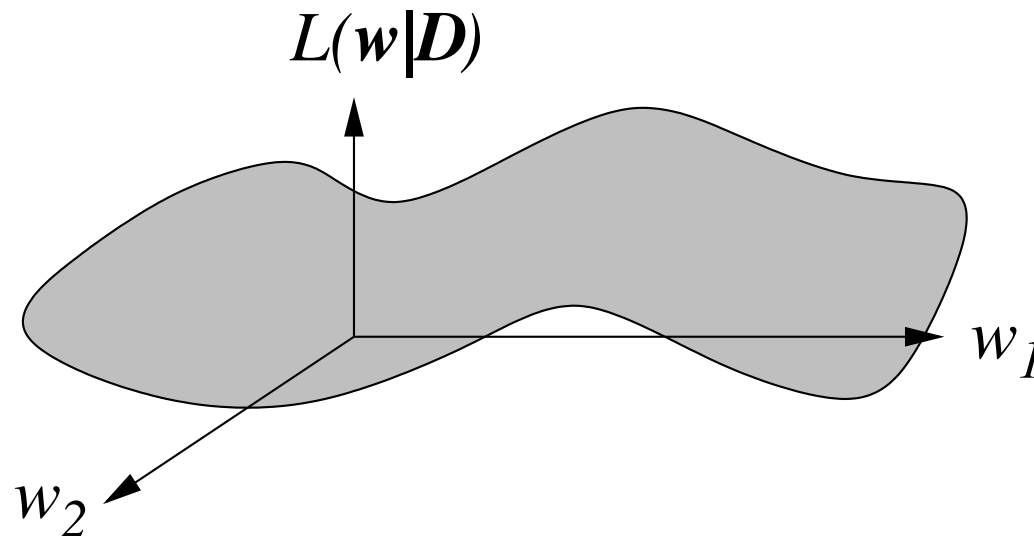
- $L(\mathbf{w}|\mathcal{D})$  is a complex function of the weights  $\mathbf{w}$



- To minimise we  $L(\mathbf{w}|\mathcal{D})$  we compute the gradient  $\nabla L(\mathbf{w}|\mathcal{D})$
- In MLP an efficient algorithm for computing the gradient is known as back-prop

# Computing Gradients

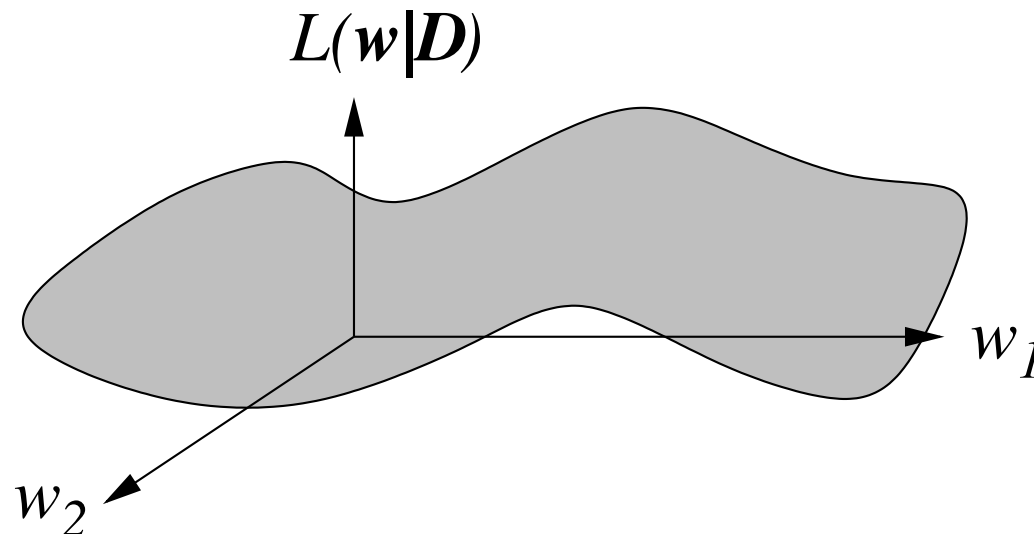
- $L(\boldsymbol{w}|\mathcal{D})$  is a complex function of the weights  $\boldsymbol{w}$



- To minimise we  $L(\boldsymbol{w}|\mathcal{D})$  we compute the gradient  $\nabla L(\boldsymbol{w}|\mathcal{D})$
- In MLP an efficient algorithm for computing the gradient is known as back-prop

# Computing Gradients

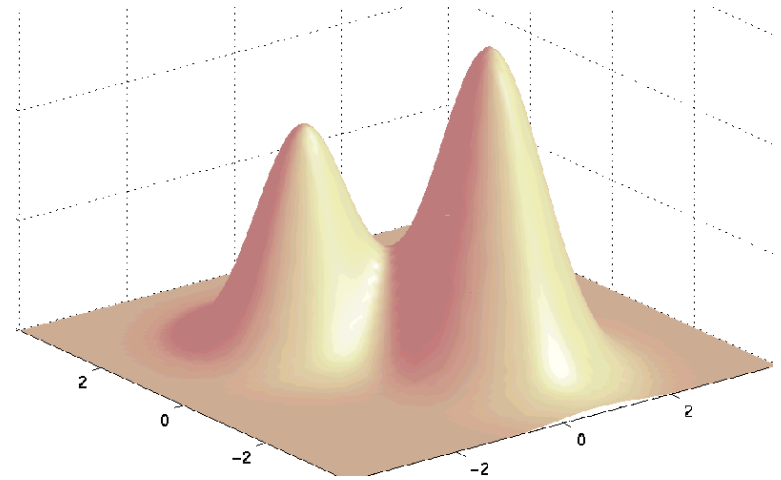
- $L(\boldsymbol{w}|\mathcal{D})$  is a complex function of the weights  $\boldsymbol{w}$



- To minimise we  $L(\boldsymbol{w}|\mathcal{D})$  we compute the gradient  $\nabla L(\boldsymbol{w}|\mathcal{D})$
- In MLP an efficient algorithm for computing the gradient is known as back-prop

# Outline

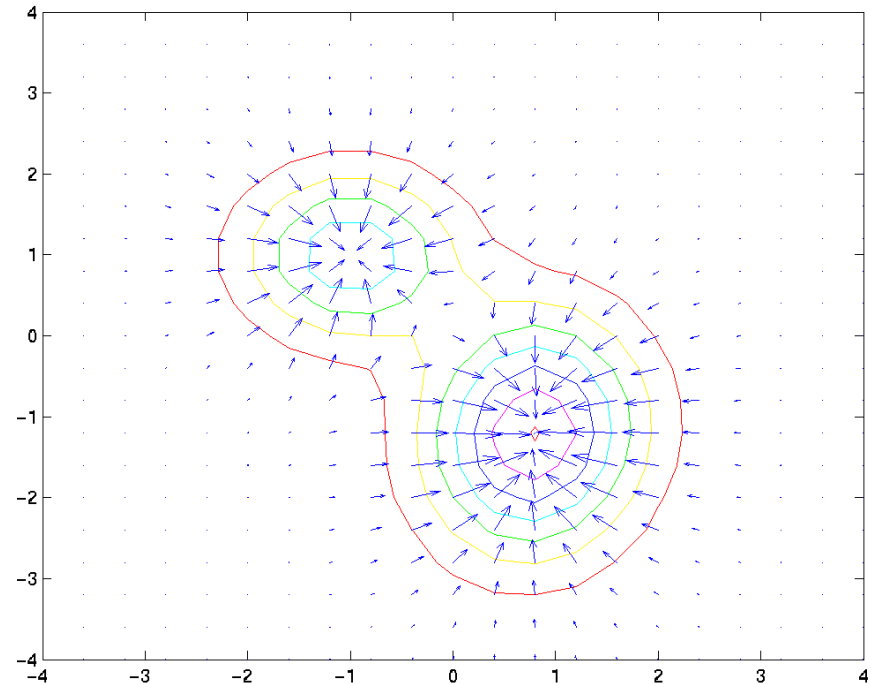
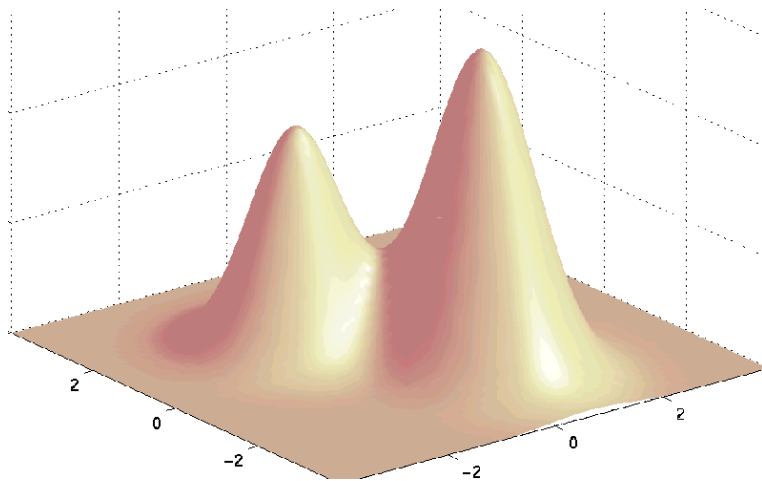
1. Motivation
2. **Gradient Descent**
3. Why Gradient Descent is Difficult



# Gradient Optimisation

- A maximum or minimum occurs when  $\nabla L(w|\mathcal{D}) = \mathbf{0}$
- E.g.

$$z = e^{-(x+1)^2-(y-1)^2} + 0.6 e^{-(x-1)^2-0.5(y+1)^2+0.1(x-3)(y-3)}$$

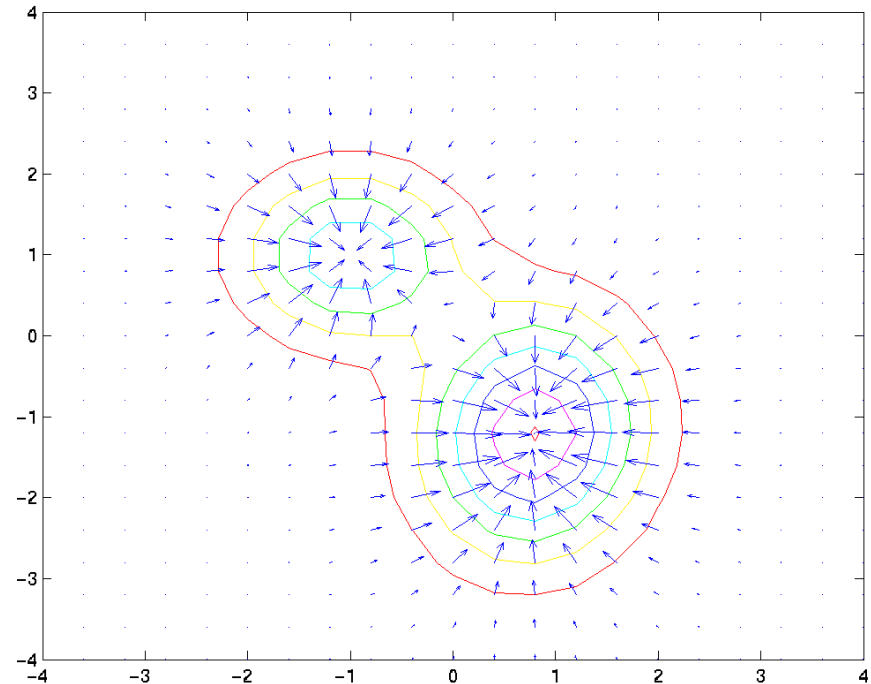
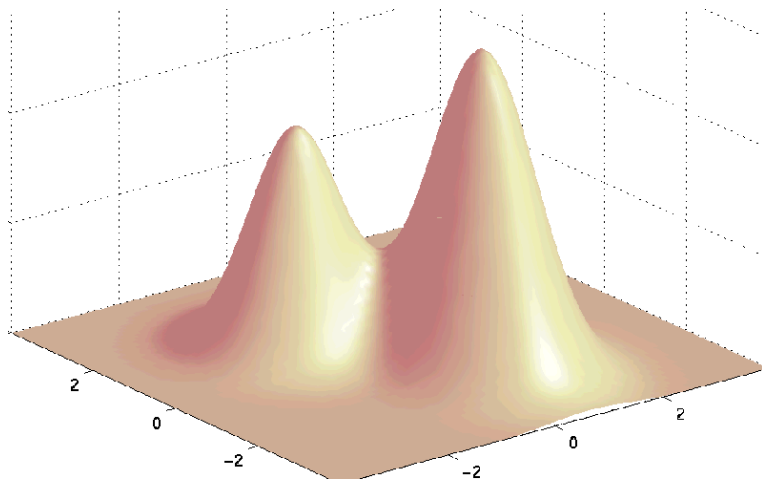




# Gradient Optimisation

- A maximum or minimum occurs when  $\nabla L(w|\mathcal{D}) = \mathbf{0}$
- E.g.

$$z = e^{-(x+1)^2-(y-1)^2} + 0.6 e^{-(x-1)^2-0.5(y+1)^2+0.1(x-3)(y-3)}$$



# Gradient Descent

- For a simple function  $L(\boldsymbol{w}|\mathcal{D})$  we can solve  $\nabla L(\boldsymbol{w}|\mathcal{D}) = \mathbf{0}$  explicitly. E.g. the linear perceptron
- For a non-linear functions we usually can't solve this set of simultaneous equations
- We can find a maximum or minimum **iteratively**
- If we know the gradient then we can follow the gradient
  - ★ Maximisation:  $\boldsymbol{w} \rightarrow \boldsymbol{w}' = \boldsymbol{w} + r \nabla L(\boldsymbol{w}|\mathcal{D})$
  - ★ Minimisation:  $\boldsymbol{w} \rightarrow \boldsymbol{w}' = \boldsymbol{w} - r \nabla L(\boldsymbol{w}|\mathcal{D})$

# Gradient Descent

- For a simple function  $L(\mathbf{w}|\mathcal{D})$  we can solve  $\nabla L(\mathbf{w}|\mathcal{D}) = \mathbf{0}$  explicitly. E.g. the linear perceptron
- For a non-linear functions we usually can't solve this set of simultaneous equations
- We can find a maximum or minimum **iteratively**
- If we know the gradient then we can follow the gradient
  - ★ Maximisation:  $\mathbf{w} \rightarrow \mathbf{w}' = \mathbf{w} + r \nabla L(\mathbf{w}|\mathcal{D})$
  - ★ Minimisation:  $\mathbf{w} \rightarrow \mathbf{w}' = \mathbf{w} - r \nabla L(\mathbf{w}|\mathcal{D})$

# Gradient Descent

- For a simple function  $L(\boldsymbol{w}|\mathcal{D})$  we can solve  $\nabla L(\boldsymbol{w}|\mathcal{D}) = \mathbf{0}$  explicitly. E.g. the linear perceptron
- For a non-linear functions we usually can't solve this set of simultaneous equations
- We can find a maximum or minimum **iteratively**
- If we know the gradient then we can follow the gradient
  - ★ Maximisation:  $\boldsymbol{w} \rightarrow \boldsymbol{w}' = \boldsymbol{w} + r \nabla L(\boldsymbol{w}|\mathcal{D})$
  - ★ Minimisation:  $\boldsymbol{w} \rightarrow \boldsymbol{w}' = \boldsymbol{w} - r \nabla L(\boldsymbol{w}|\mathcal{D})$

# Gradient Descent

- For a simple function  $L(\mathbf{w}|\mathcal{D})$  we can solve  $\nabla L(\mathbf{w}|\mathcal{D}) = \mathbf{0}$  explicitly. E.g. the linear perceptron
- For a non-linear functions we usually can't solve this set of simultaneous equations
- We can find a maximum or minimum **iteratively**
- If we know the gradient then we can follow the gradient
  - ★ Maximisation:  $\mathbf{w} \rightarrow \mathbf{w}' = \mathbf{w} + r \nabla L(\mathbf{w}|\mathcal{D})$
  - ★ Minimisation:  $\mathbf{w} \rightarrow \mathbf{w}' = \mathbf{w} - r \nabla L(\mathbf{w}|\mathcal{D})$

# Gradient Descent

- For a simple function  $L(\boldsymbol{w}|\mathcal{D})$  we can solve  $\nabla L(\boldsymbol{w}|\mathcal{D}) = \mathbf{0}$  explicitly. E.g. the linear perceptron
- For a non-linear functions we usually can't solve this set of simultaneous equations
- We can find a maximum or minimum **iteratively**
- If we know the gradient then we can follow the gradient
  - ★ Maximisation:  $\boldsymbol{w} \rightarrow \boldsymbol{w}' = \boldsymbol{w} + r \nabla L(\boldsymbol{w}|\mathcal{D})$
  - ★ Minimisation:  $\boldsymbol{w} \rightarrow \boldsymbol{w}' = \boldsymbol{w} - r \nabla L(\boldsymbol{w}|\mathcal{D})$

# Gradient Descent

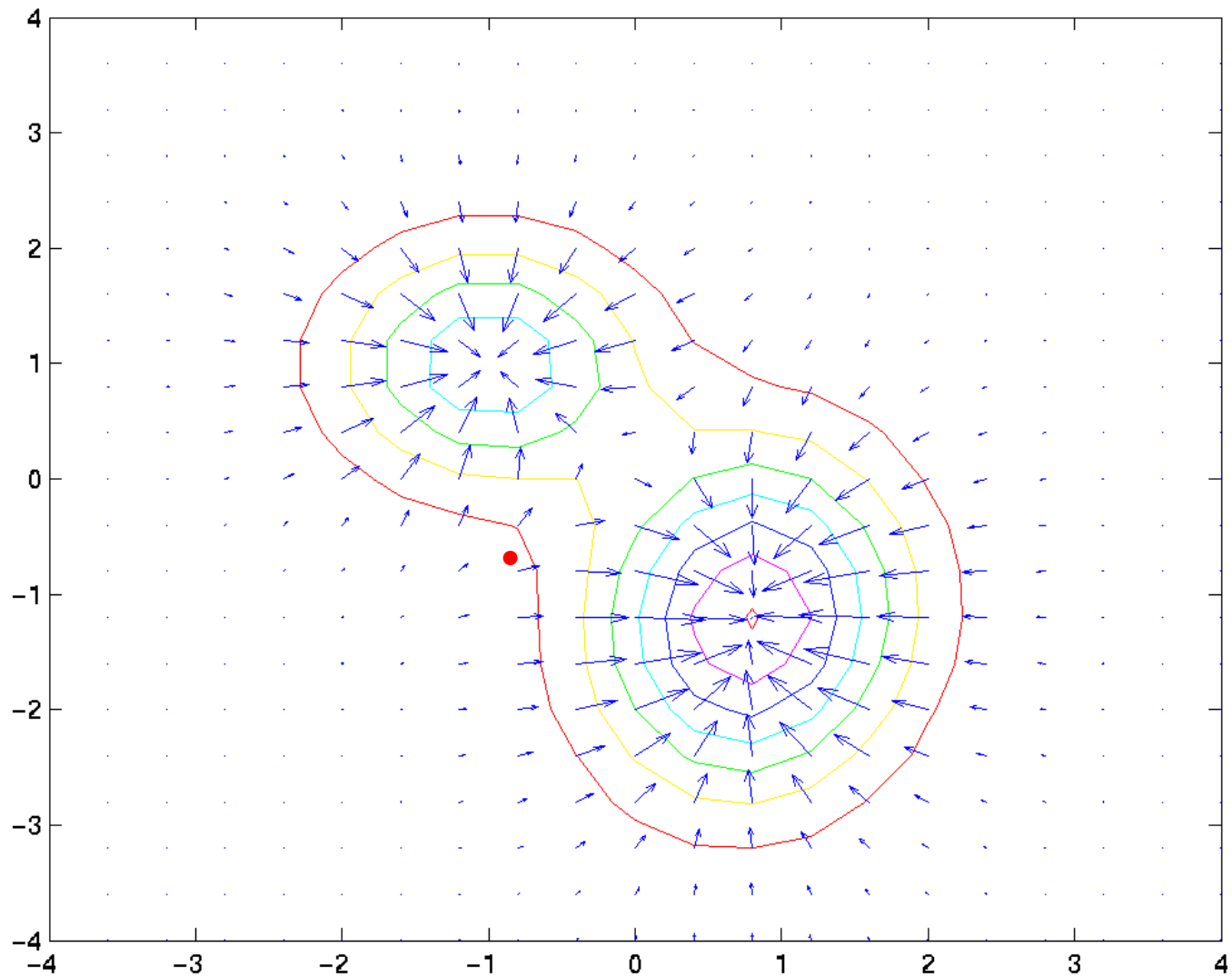
- For a simple function  $L(\boldsymbol{w}|\mathcal{D})$  we can solve  $\nabla L(\boldsymbol{w}|\mathcal{D}) = \mathbf{0}$  explicitly. E.g. the linear perceptron
- For a non-linear functions we usually can't solve this set of simultaneous equations
- We can find a maximum or minimum **iteratively**
- If we know the gradient then we can follow the gradient
  - ★ Maximisation:  $\boldsymbol{w} \rightarrow \boldsymbol{w}' = \boldsymbol{w} + r \nabla L(\boldsymbol{w}|\mathcal{D})$
  - ★ Minimisation:  $\boldsymbol{w} \rightarrow \boldsymbol{w}' = \boldsymbol{w} - r \nabla L(\boldsymbol{w}|\mathcal{D})$

# Gradient Descent

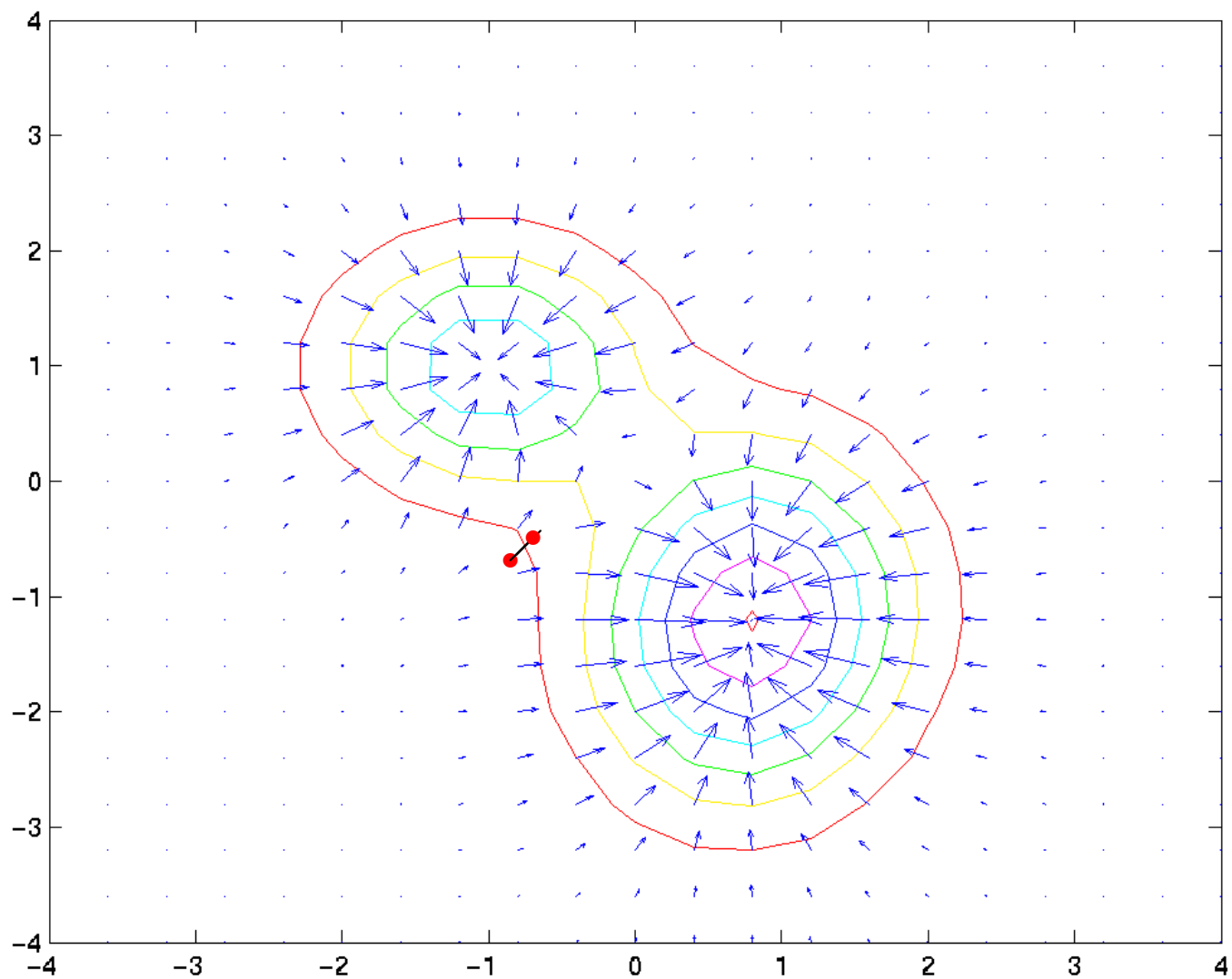
- For a simple function  $L(\boldsymbol{w}|\mathcal{D})$  we can solve  $\nabla L(\boldsymbol{w}|\mathcal{D}) = \mathbf{0}$  explicitly. E.g. the linear perceptron
- For a non-linear functions we usually can't solve this set of simultaneous equations
- We can find a maximum or minimum **iteratively**
- If we know the gradient then we can follow the gradient
  - ★ Maximisation:  $\boldsymbol{w} \rightarrow \boldsymbol{w}' = \boldsymbol{w} + r \nabla L(\boldsymbol{w}|\mathcal{D})$
  - ★ Minimisation:  $\boldsymbol{w} \rightarrow \boldsymbol{w}' = \boldsymbol{w} - r \nabla L(\boldsymbol{w}|\mathcal{D})$



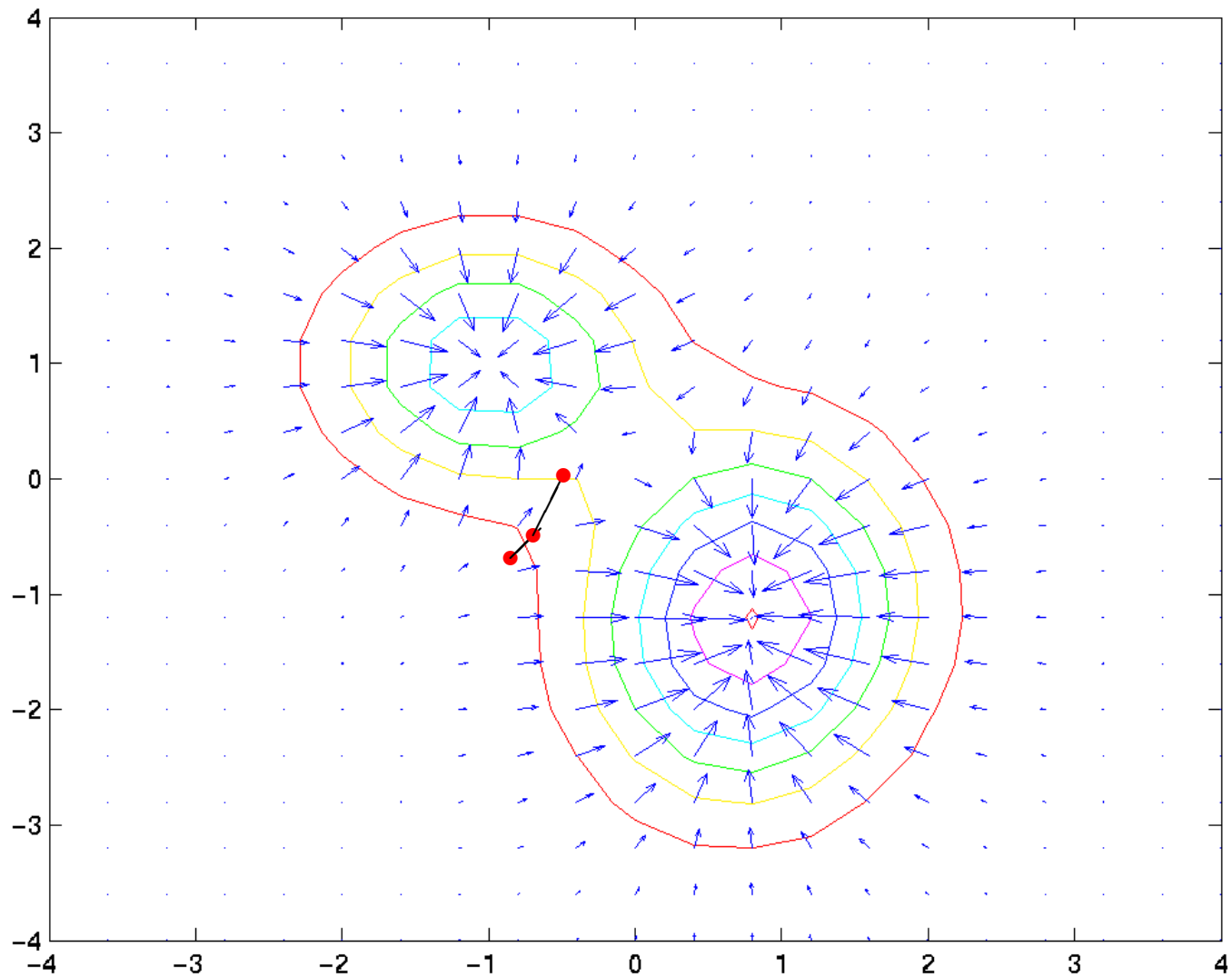
# Hill-Climbing



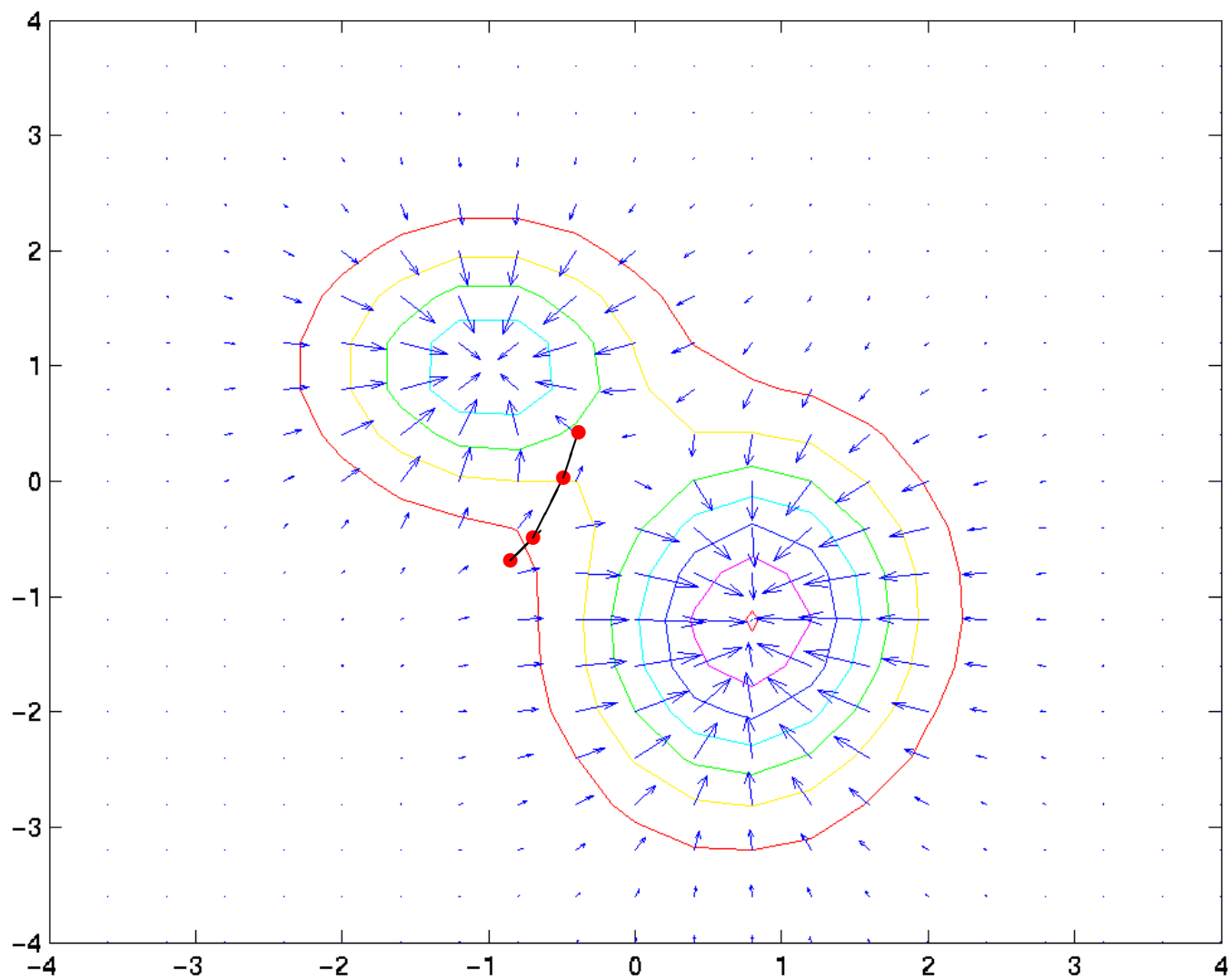
# Hill-Climbing



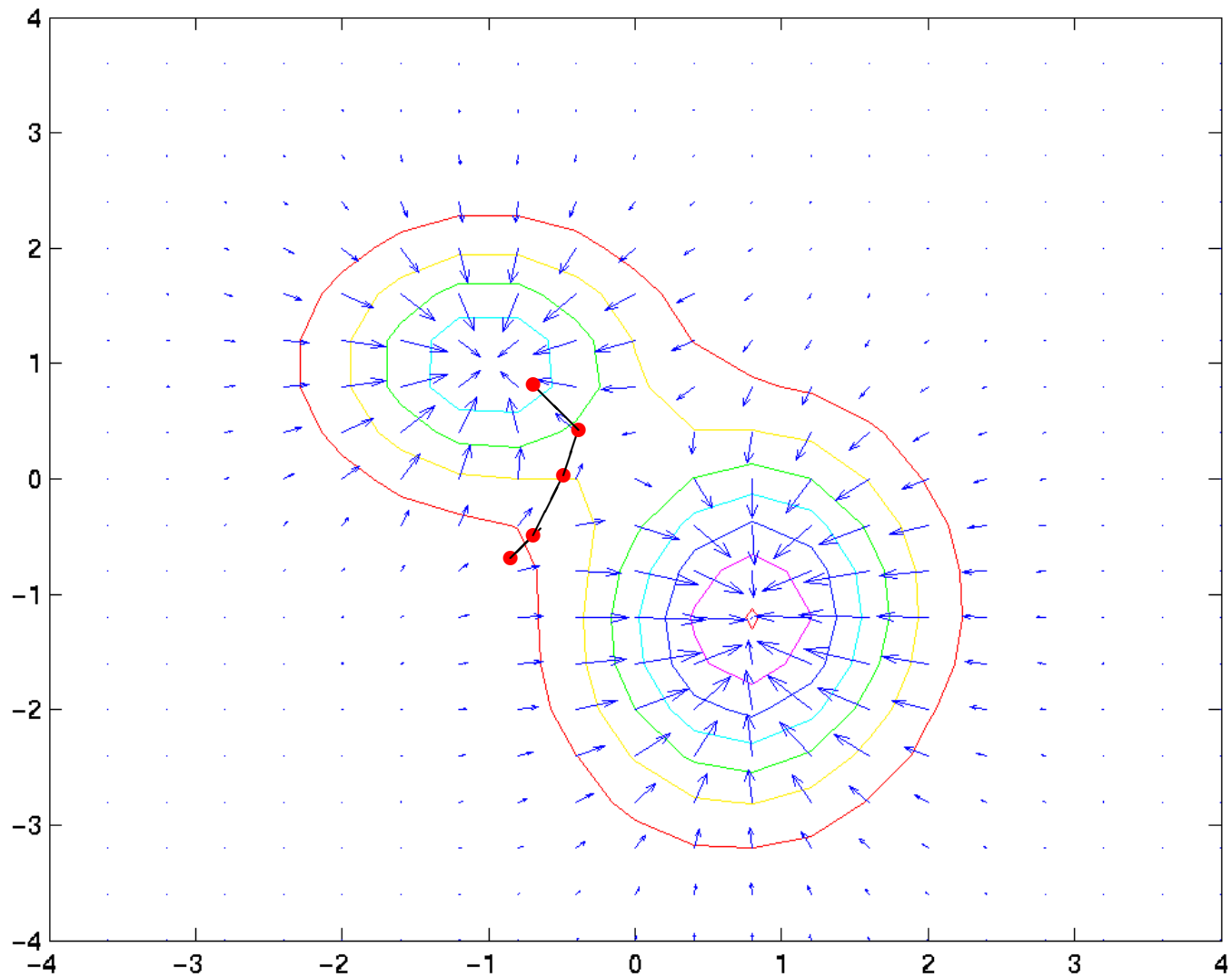
# Hill-Climbing



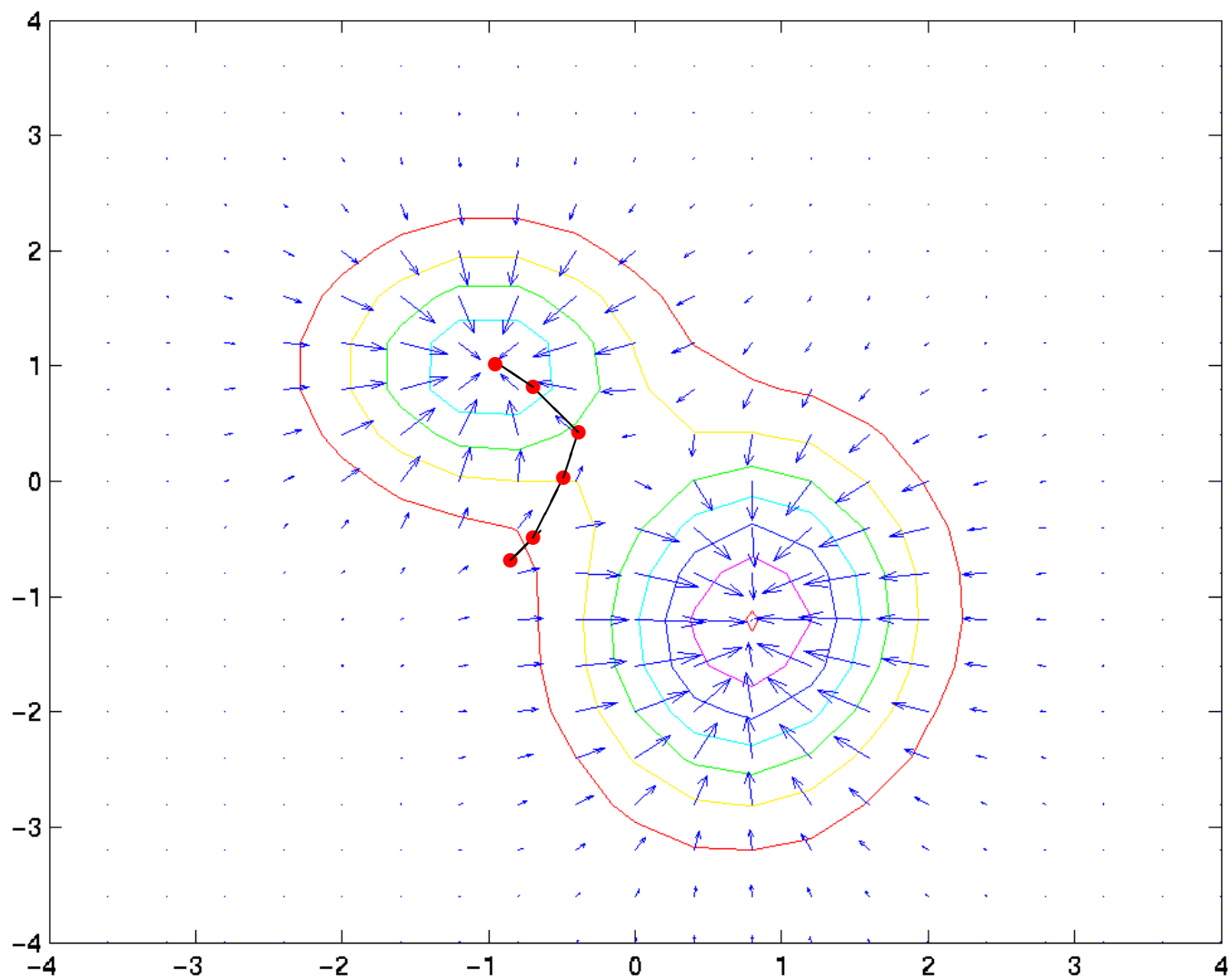
# Hill-Climbing



# Hill-Climbing

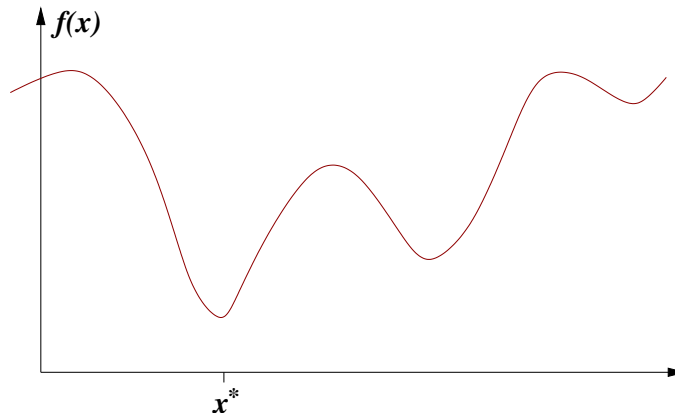


# Hill-Climbing



# What Goes Right

- Almost all minima are quadratic (Morse's theorem)



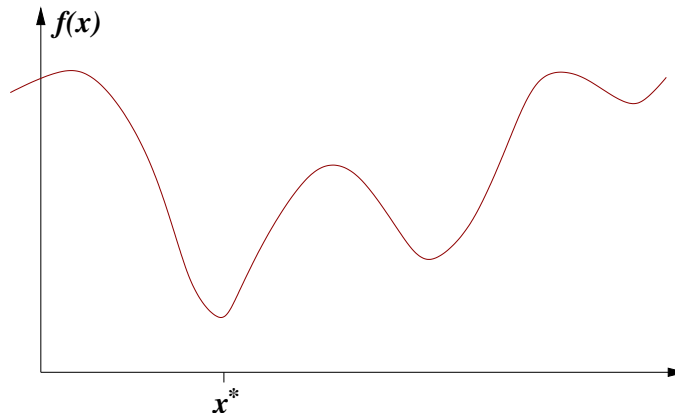
- Taylor expanding around a minimum  $x^*$

$$\begin{aligned} f(x) &= f(x^*) + (x - x^*) f'(x^*) + \frac{1}{2} (x - x^*)^2 f''(x^*) + \dots \\ &= f(x^*) + \frac{1}{2} (x - x^*)^2 f''(x^*) + \frac{1}{3!} (x - x^*)^3 f'''(x^*) + \dots \end{aligned}$$

- If  $x - x^*$  is sufficiently small the higher order terms are negligible

# What Goes Right

- Almost all minima are quadratic (Morse's theorem)



- Taylor expanding around a minimum  $x^*$

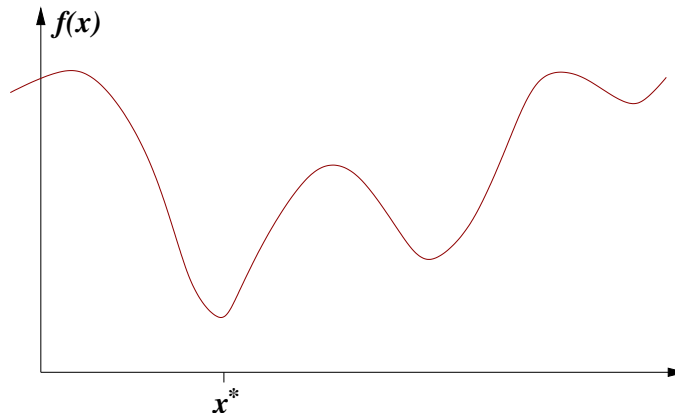
$$\begin{aligned} f(x) &= f(x^*) + (x - x^*) f'(x^*) + \frac{1}{2} (x - x^*)^2 f''(x^*) + \dots \\ &= f(x^*) + \frac{1}{2} (x - x^*)^2 f''(x^*) + \frac{1}{3!} (x - x^*)^3 f'''(x^*) + \dots \end{aligned}$$

- If  $x - x^*$  is sufficiently small the higher order terms are negligible



# What Goes Right

- Almost all minima are quadratic (Morse's theorem)



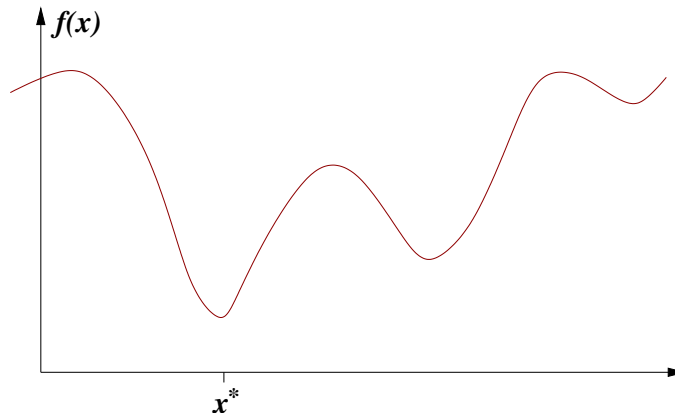
- Taylor expanding around a minimum  $x^*$

$$\begin{aligned} f(x) &= f(x^*) + (x - x^*) f'(x^*) + \frac{1}{2} (x - x^*)^2 f''(x^*) + \dots \\ &= f(x^*) + \frac{1}{2} (x - x^*)^2 f''(x^*) + \frac{1}{3!} (x - x^*)^3 f'''(x^*) + \dots \end{aligned}$$

- If  $x - x^*$  is sufficiently small the higher order terms are negligible

# What Goes Right

- Almost all minima are quadratic (Morse's theorem)



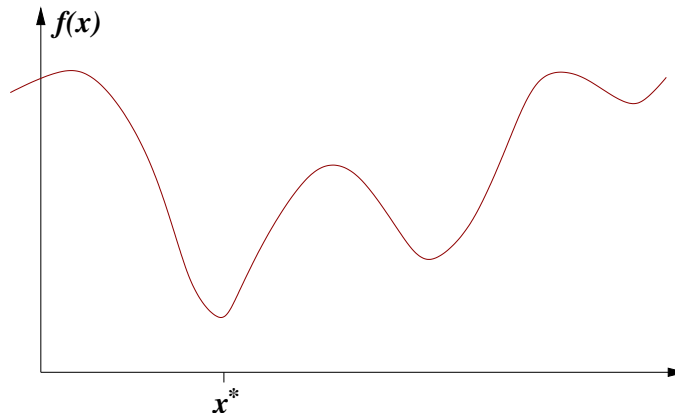
- Taylor expanding around a minimum  $x^*$

$$\begin{aligned} f(x) &= f(x^*) + (x - x^*) f'(x^*) + \frac{1}{2} (x - x^*)^2 f''(x^*) + \dots \\ &= f(x^*) + \frac{1}{2} (x - x^*)^2 f''(x^*) + \frac{1}{3!} (x - x^*)^3 f'''(x^*) + \dots \end{aligned}$$

- If  $x - x^*$  is sufficiently small the higher order terms are negligible

# What Goes Right

- Almost all minima are quadratic (Morse's theorem)



- Taylor expanding around a minimum  $x^*$

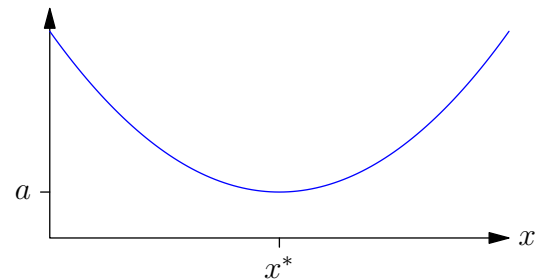
$$\begin{aligned} f(x) &= f(x^*) + (x - x^*) f'(x^*) + \frac{1}{2} (x - x^*)^2 f''(x^*) + \dots \\ &= f(x^*) + \frac{1}{2} (x - x^*)^2 f''(x^*) + \frac{1}{3!} (x - x^*)^3 f'''(x^*) + \dots \end{aligned}$$

- If  $x - x^*$  is sufficiently small the higher order terms are negligible

# Newton's Method

- If we were in a quadratic minimum

$$f(x) = a + \frac{b}{2} (x - x^*)^2$$



- then

$$f'(x) = b (x - x^*), \quad f''(x) = b$$

- so

$$x - x^* = \frac{f'(x)}{b} = \frac{f'(x)}{f''(x)}$$

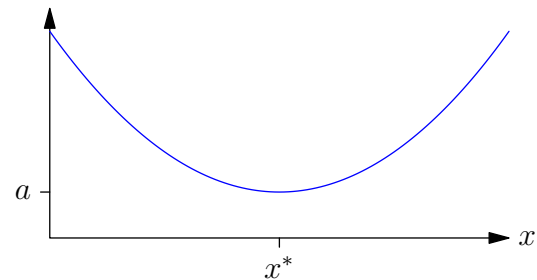
- or

$$x^* = x - \frac{f'(x)}{f''(x)}$$

# Newton's Method

- If we were in a quadratic minimum

$$f(x) = a + \frac{b}{2} (x - x^*)^2$$



- then

$$f'(x) = b (x - x^*), \quad f''(x) = b$$

- so

$$x - x^* = \frac{f'(x)}{b} = \frac{f'(x)}{f''(x)}$$

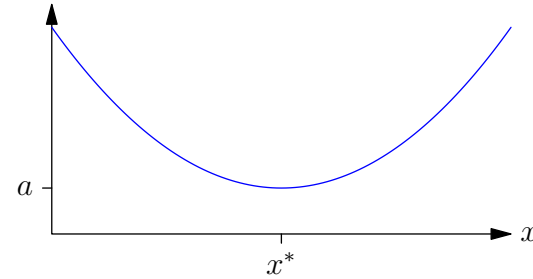
- or

$$x^* = x - \frac{f'(x)}{f''(x)}$$

# Newton's Method

- If we were in a quadratic minimum

$$f(x) = a + \frac{b}{2} (x - x^*)^2$$



- then

$$f'(x) = b (x - x^*), \quad f''(x) = b$$

- so

$$x - x^* = \frac{f'(x)}{b} = \frac{f'(x)}{f''(x)}$$

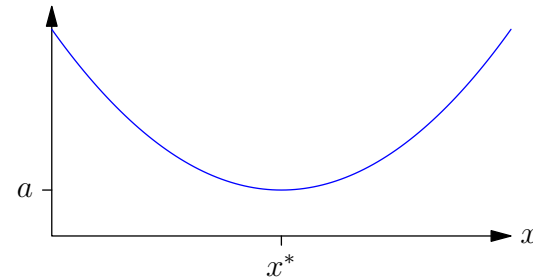
- or

$$x^* = x - \frac{f'(x)}{f''(x)}$$

# Newton's Method

- If we were in a quadratic minimum

$$f(x) = a + \frac{b}{2} (x - x^*)^2$$



- then

$$f'(x) = b (x - x^*), \quad f''(x) = b$$

- so

$$x - x^* = \frac{f'(x)}{b} = \frac{f'(x)}{f''(x)}$$

- or

$$x^* = x - \frac{f'(x)}{f''(x)}$$

# Newton's Method

- This is Newton's methods
- For non-quadratic functions Newton's method converges **quadratically** provided we are sufficiently close to a minimum
- If we are at a distance  $x - x^* = \epsilon$  from the minima then after one cycle we will be a distance  $\epsilon^2$  after two cycles we will be at a distance  $\epsilon^4$ , etc.
- If we are too far from a minimum we might go anywhere!
- We should follow the gradient until we are near the minimum



# Newton's Method

- This is Newton's methods
- For non-quadratic functions Newtons method converges **quadratically** provided we are sufficiently close to a minimum
- If we are at a distance  $x - x^* = \epsilon$  from the minima then after one cycle we will be a distance  $\epsilon^2$  after two cycles we will be at a distance  $\epsilon^4$ , etc.
- If we are too far from a minimum we might go anywhere!
- We should follow the gradient until we are near the minimum

# Newton's Method

- This is Newton's methods
- For non-quadratic functions Newtons method converges **quadratically** provided we are sufficiently close to a minimum
- If we are at a distance  $x - x^* = \epsilon$  from the minima then after one cycle we will be a distance  $\epsilon^2$  after two cycles we will be at a distance  $\epsilon^4$ , etc.
- If we are too far from a minimum we might go anywhere!
- We should follow the gradient until we are near the minimum

# Newton's Method

- This is Newton's methods
- For non-quadratic functions Newtons method converges **quadratically** provided we are sufficiently close to a minimum
- If we are at a distance  $x - x^* = \epsilon$  from the minima then after one cycle we will be a distance  $\epsilon^2$  **after two cycles we will be at a distance  $\epsilon^4$ , etc.**
- If we are too far from a minimum we might go anywhere!
- We should follow the gradient until we are near the minimum

# Newton's Method

- This is Newton's methods
- For non-quadratic functions Newtons method converges **quadratically** provided we are sufficiently close to a minimum
- If we are at a distance  $x - x^* = \epsilon$  from the minima then after one cycle we will be a distance  $\epsilon^2$  after two cycles we will be at a distance  $\epsilon^4$ , etc.
- If we are too far from a minimum we might go anywhere!
- We should follow the gradient until we are near the minimum

# Newton's Method

- This is Newton's methods
- For non-quadratic functions Newtons method converges **quadratically** provided we are sufficiently close to a minimum
- If we are at a distance  $x - x^* = \epsilon$  from the minima then after one cycle we will be a distance  $\epsilon^2$  after two cycles we will be at a distance  $\epsilon^4$ , etc.
- If we are too far from a minimum we might go anywhere!
- We should follow the gradient until we are near the minimum

# Taylor's Expansion in High Dimensions

- We can generalise these results to many dimensions
- The Taylor expansion of a function  $f(\mathbf{x})$  about  $\mathbf{x}_0$

$$f(\mathbf{x}) = f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^\top \nabla f(\mathbf{x}_0) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_0)^\top \mathbf{H} (\mathbf{x} - \mathbf{x}_0) + \dots$$

where  $\mathbf{H}$  is the **Hessian** matrix with elements

$$H_{ij} = \frac{\partial^2 f(\mathbf{x}_0)}{\partial x_i \partial x_j}$$

- Newton's method in high dimension is

$$\mathbf{x}^* = \mathbf{x} - \mathbf{H}^{-1} \nabla f(\mathbf{x})$$

# Taylor's Expansion in High Dimensions

- We can generalise these results to many dimensions
- The Taylor expansion of a function  $f(\mathbf{x})$  about  $\mathbf{x}_0$

$$f(\mathbf{x}) = f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^\top \nabla f(\mathbf{x}_0) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_0)^\top \mathbf{H} (\mathbf{x} - \mathbf{x}_0) + \dots$$

where  $\mathbf{H}$  is the **Hessian** matrix with elements

$$H_{ij} = \frac{\partial^2 f(\mathbf{x}_0)}{\partial x_i \partial x_j}$$

- Newton's method in high dimension is

$$\mathbf{x}^* = \mathbf{x} - \mathbf{H}^{-1} \nabla f(\mathbf{x})$$

# Taylor's Expansion in High Dimensions

- We can generalise these results to many dimensions
- The Taylor expansion of a function  $f(\mathbf{x})$  about  $\mathbf{x}_0$

$$f(\mathbf{x}) = f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^\top \nabla f(\mathbf{x}_0) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_0)^\top \mathbf{H} (\mathbf{x} - \mathbf{x}_0) + \dots$$

where  $\mathbf{H}$  is the **Hessian** matrix with elements

$$H_{ij} = \frac{\partial^2 f(\mathbf{x}_0)}{\partial x_i \partial x_j}$$

- Newton's method in high dimension is

$$\mathbf{x}^* = \mathbf{x} - \mathbf{H}^{-1} \nabla f(\mathbf{x})$$



# Taylor's Expansion in High Dimensions

- We can generalise these results to many dimensions
- The Taylor expansion of a function  $f(\mathbf{x})$  about  $\mathbf{x}_0$

$$f(\mathbf{x}) = f(\mathbf{x}_0) + (\mathbf{x} - \mathbf{x}_0)^\top \nabla f(\mathbf{x}_0) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_0)^\top \mathbf{H} (\mathbf{x} - \mathbf{x}_0) + \dots$$

where  $\mathbf{H}$  is the **Hessian** matrix with elements

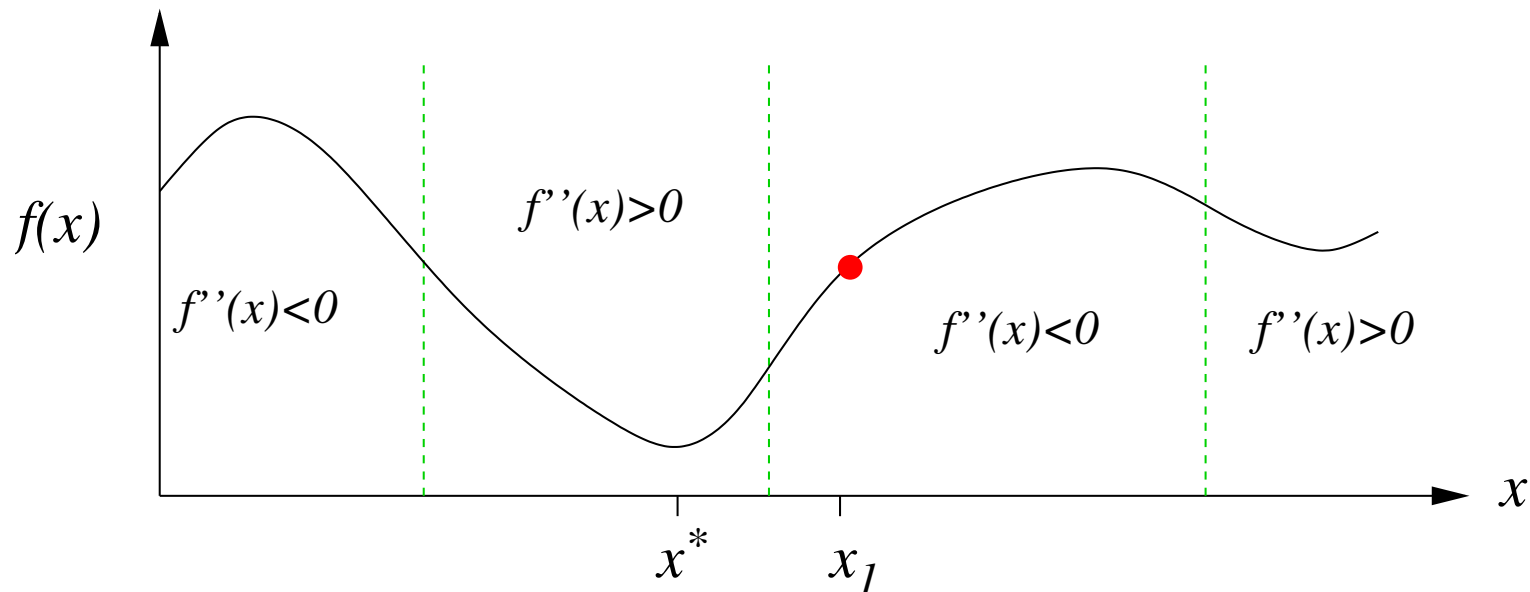
$$H_{ij} = \frac{\partial^2 f(\mathbf{x}_0)}{\partial x_i \partial x_j}$$

- Newton's method in high dimension is

$$\mathbf{x}^* = \mathbf{x} - \mathbf{H}^{-1} \nabla f(\mathbf{x})$$

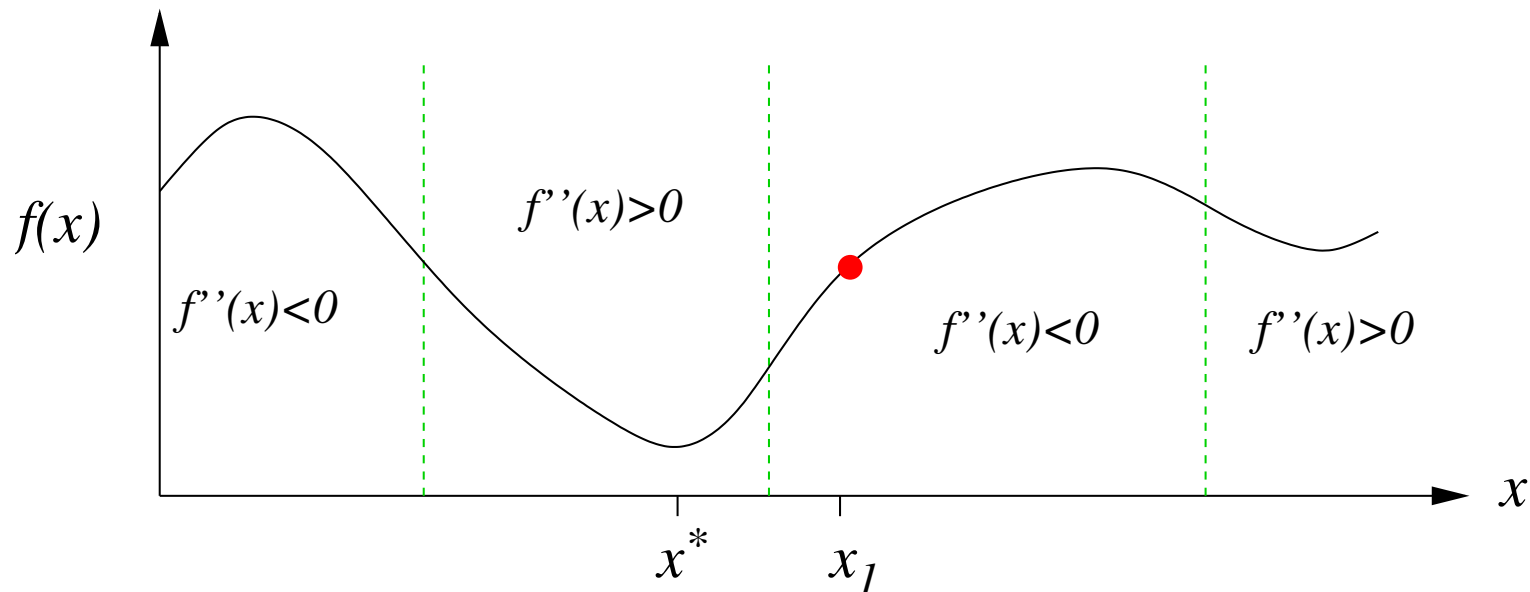
# Using the Second Derivative

- If we are optimising  $N$  parameters the Hessian is an  $N \times N$  matrix
- It is time-consuming to compute (and prone to errors when coding)
- Away from minima they can be misleading



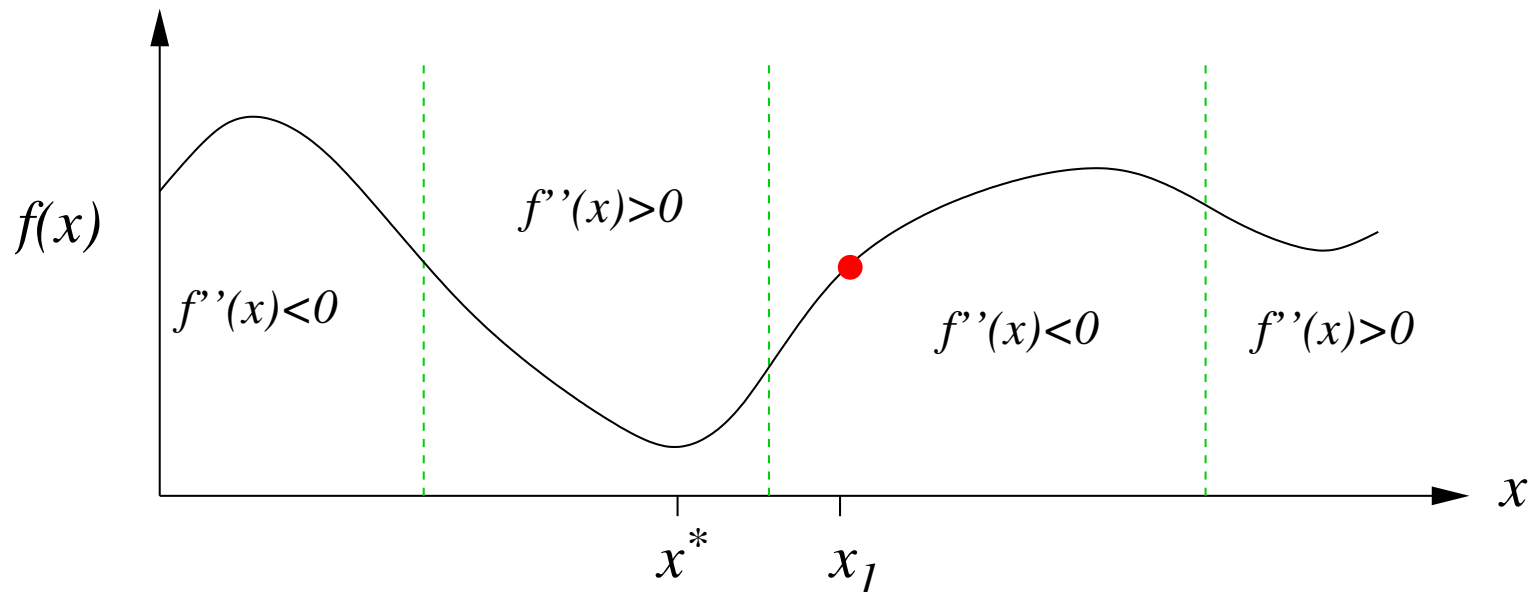
# Using the Second Derivative

- If we are optimising  $N$  parameters the Hessian is an  $N \times N$  matrix
- It is time-consuming to compute (and prone to errors when coding)
- Away from minima they can be misleading



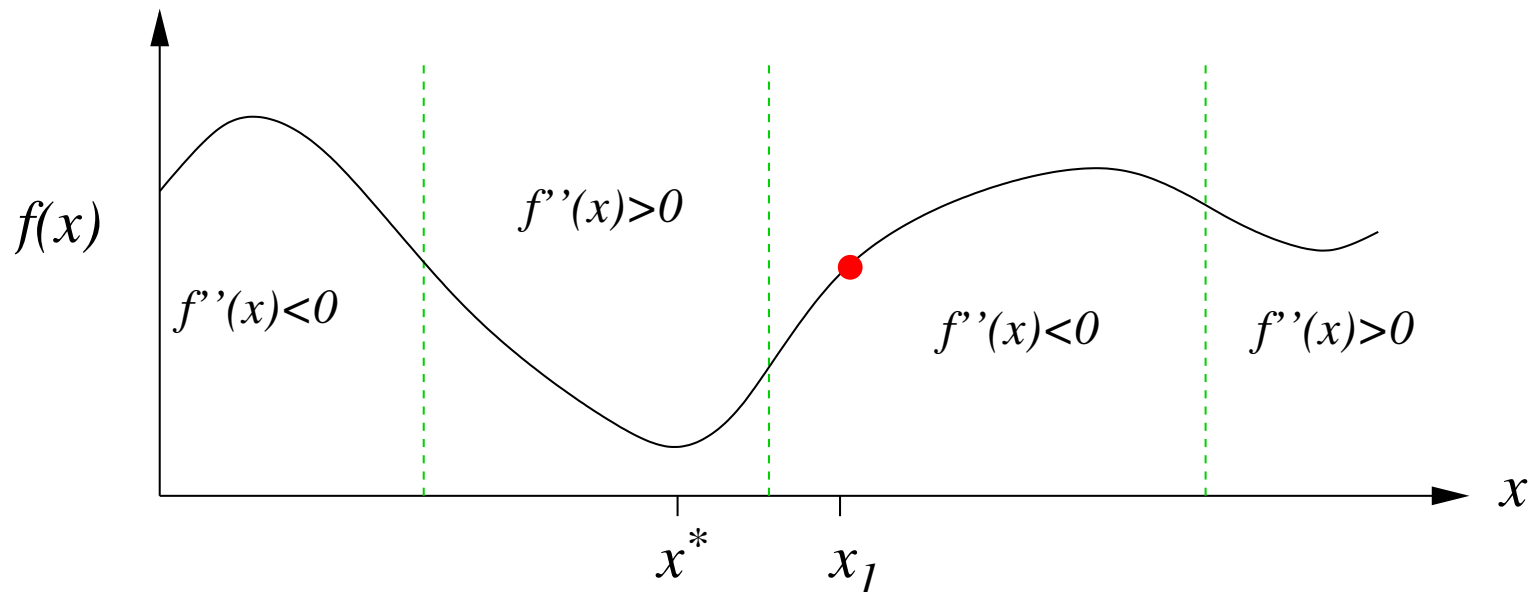
# Using the Second Derivative

- If we are optimising  $N$  parameters the Hessian is an  $N \times N$  matrix
- It is time-consuming to compute (and prone to errors when coding)
- Away from minima they can be misleading



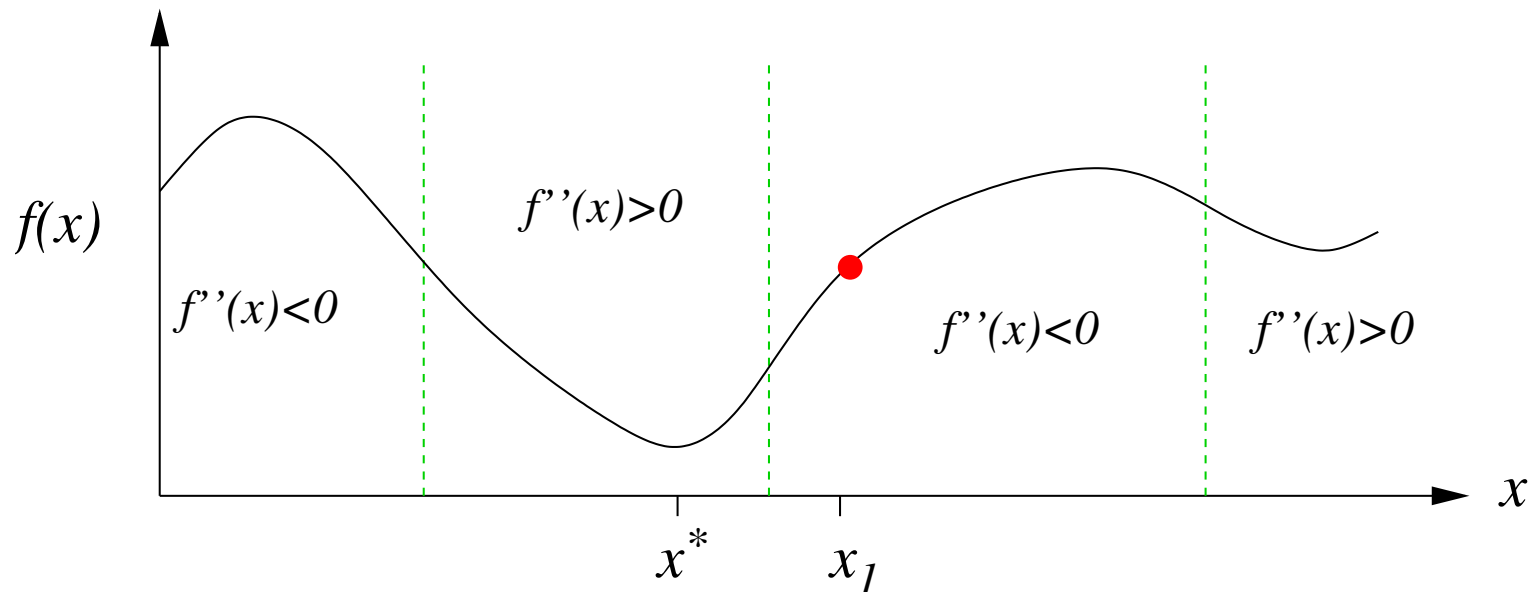
# Using the Second Derivative

- If we are optimising  $N$  parameters the Hessian is an  $N \times N$  matrix
- It is time-consuming to compute (and prone to errors when coding)—for deep learning it is impossible even to store the Hessian
- Away from minima they can be misleading



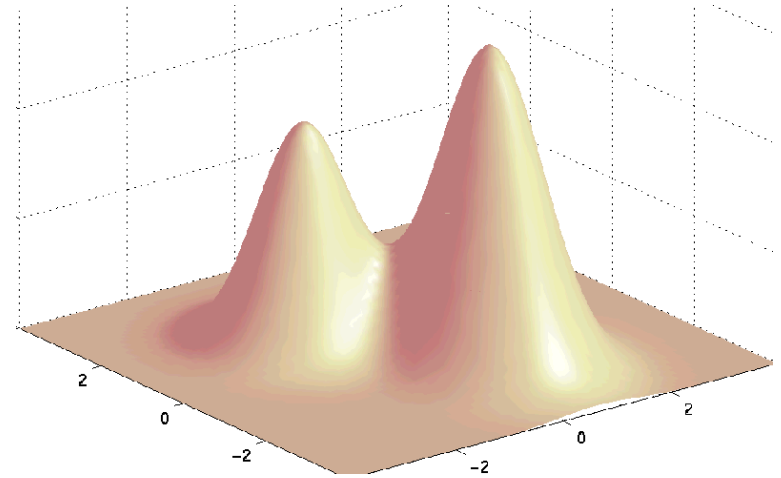
# Using the Second Derivative

- If we are optimising  $N$  parameters the Hessian is an  $N \times N$  matrix
- It is time-consuming to compute (and prone to errors when coding)—for deep learning it is impossible even to store the Hessian
- Away from minima they can be misleading



# Outline

1. Motivation
2. Gradient Descent
3. **Why Gradient Descent is Difficult**



# Step Size

- Gradient descent

$$\boldsymbol{x}' = \boldsymbol{x} - r \nabla f(\boldsymbol{x})$$

- Need to choose the learning rate of step size,  $r$
- Too small steps takes lots of time
- Too large steps takes you away from a minimum



# Step Size

- Gradient descent

$$\boldsymbol{x}' = \boldsymbol{x} - r \nabla f(\boldsymbol{x})$$

- Need to choose the learning rate of step size,  $r$
- Too small steps takes lots of time
- Too large steps takes you away from a minimum

# Step Size

- Gradient descent

$$\boldsymbol{x}' = \boldsymbol{x} - r \nabla f(\boldsymbol{x})$$

- Need to choose the learning rate of step size,  $r$
- Too small steps takes lots of time
- Too large steps takes you away from a minimum

# Step Size

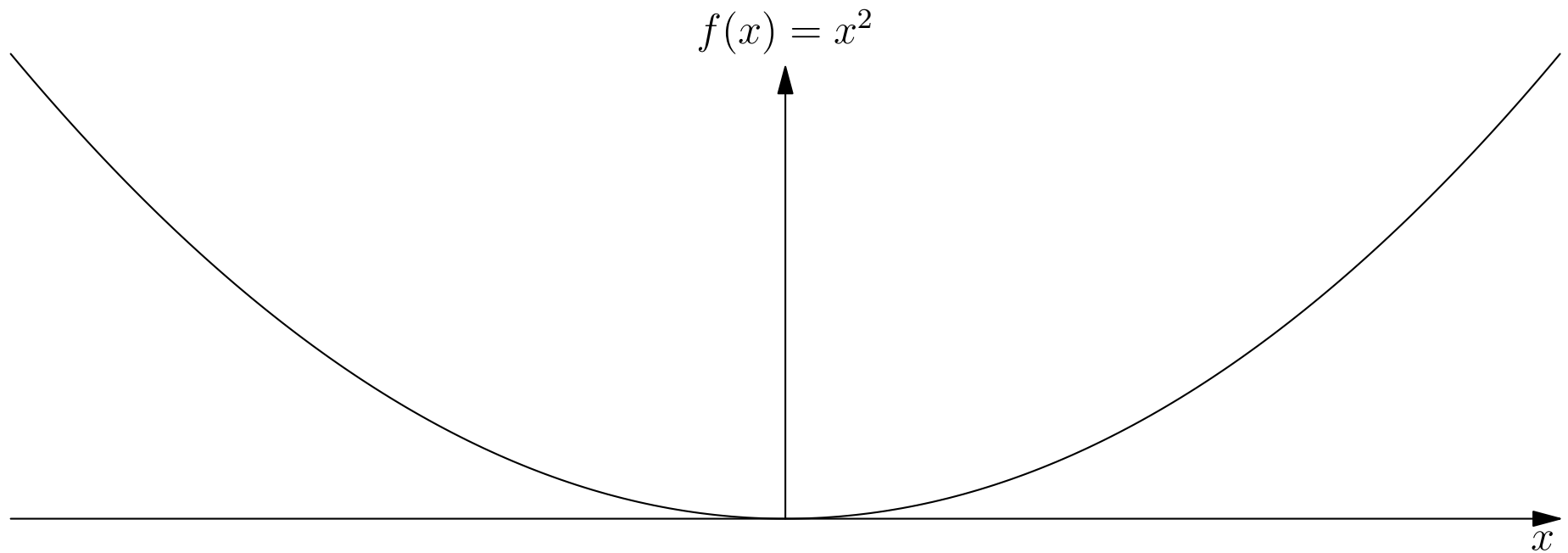
- Gradient descent

$$\mathbf{x}' = \mathbf{x} - r \nabla f(\mathbf{x})$$

- Need to choose the learning rate of step size,  $r$
- Too small steps takes lots of time
- Too large steps takes you away from a minimum

# Step Size

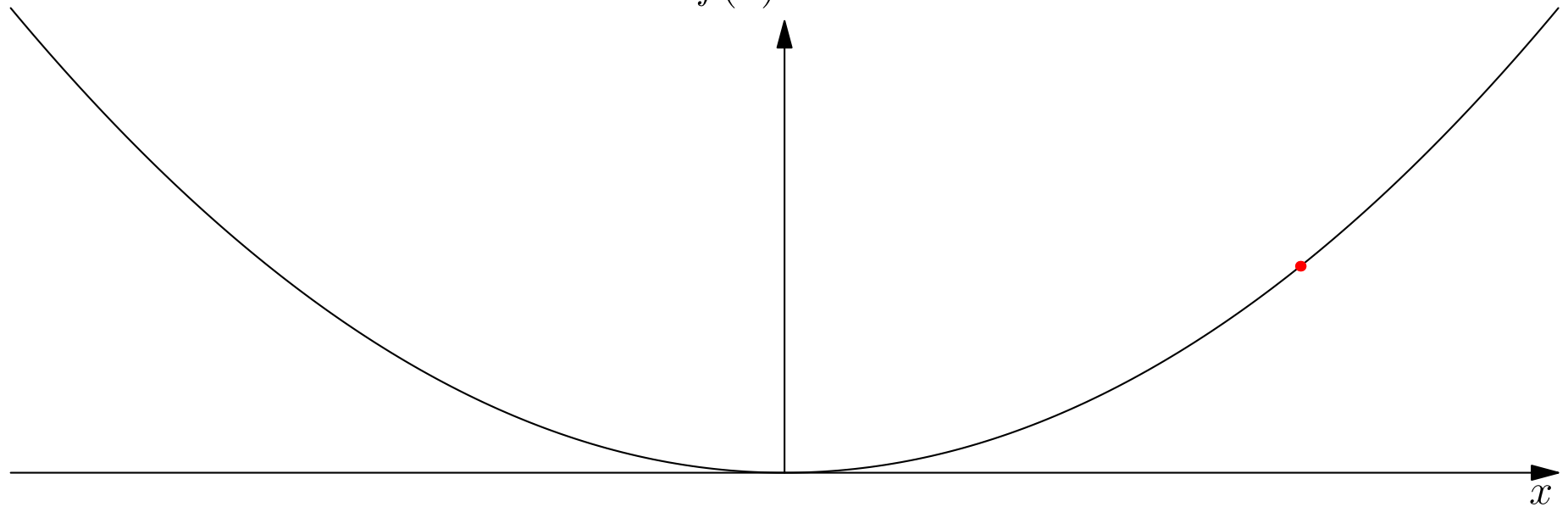
$$x \leftarrow x - r f'(x)$$



# Step Size

$$x \leftarrow x - r f'(x)$$

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

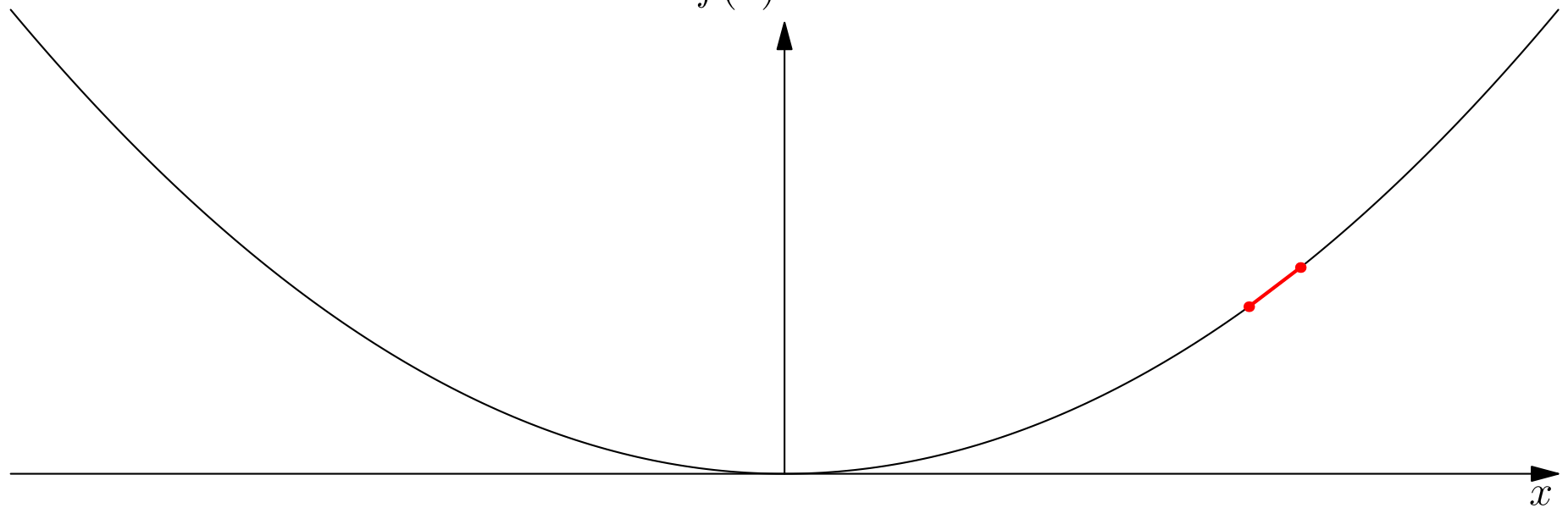


$$r = 0.05$$

# Step Size

$$x \leftarrow x - r f'(x)$$

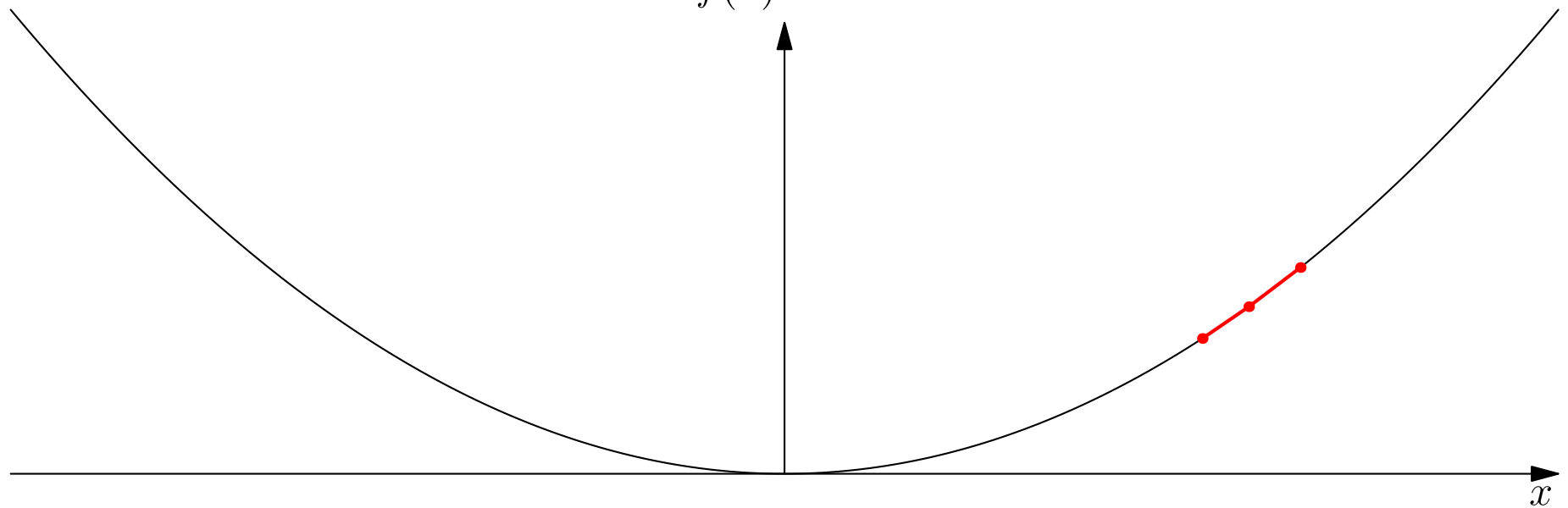
$$f(x) = x^2$$



# Step Size

$$x \leftarrow x - r f'(x)$$

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

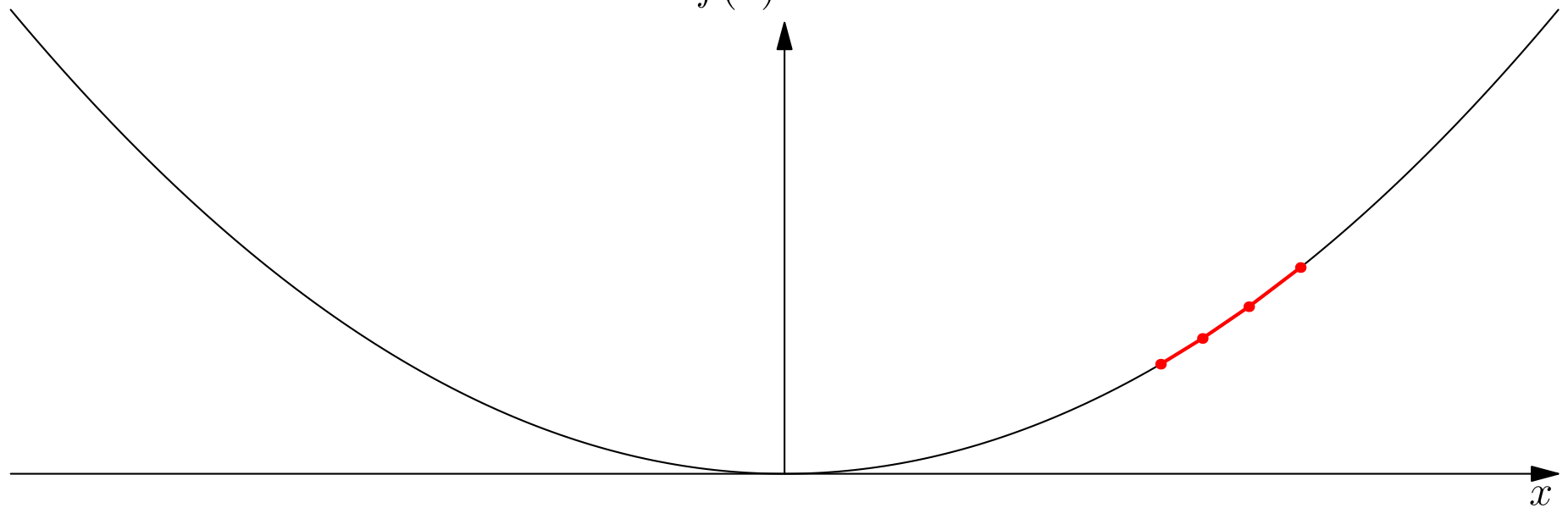


$$r = 0.05$$

# Step Size

$$x \leftarrow x - r f'(x)$$

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



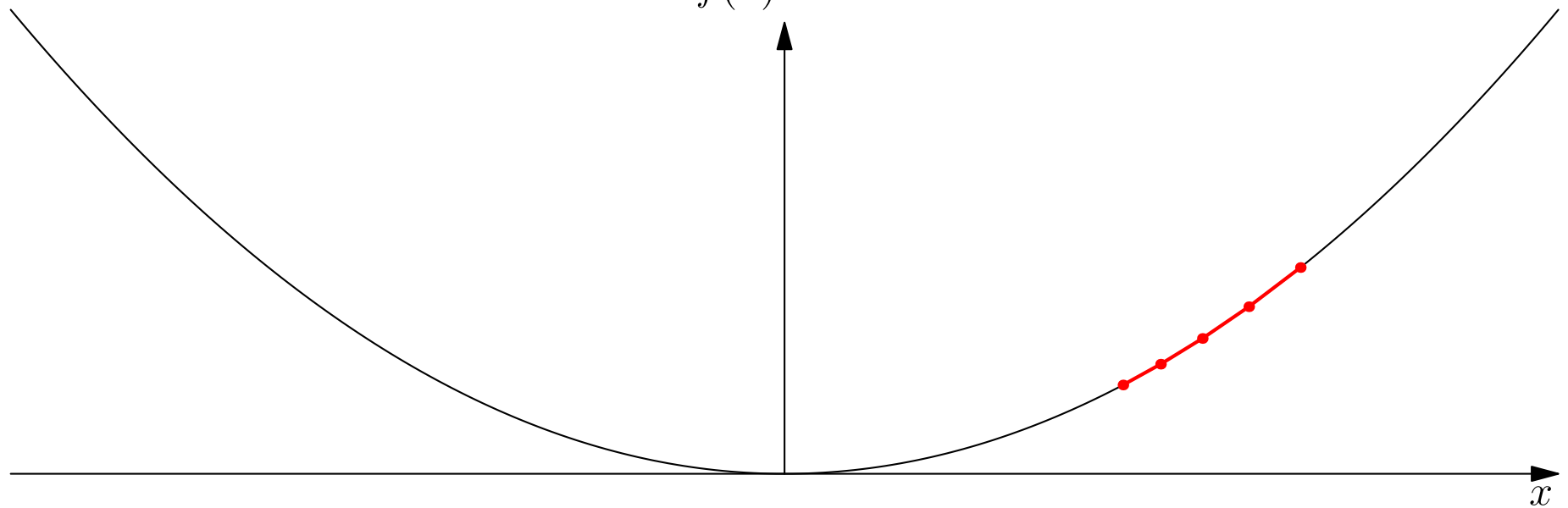
$$r = 0.05$$



# Step Size

$$x \leftarrow x - r f'(x)$$

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

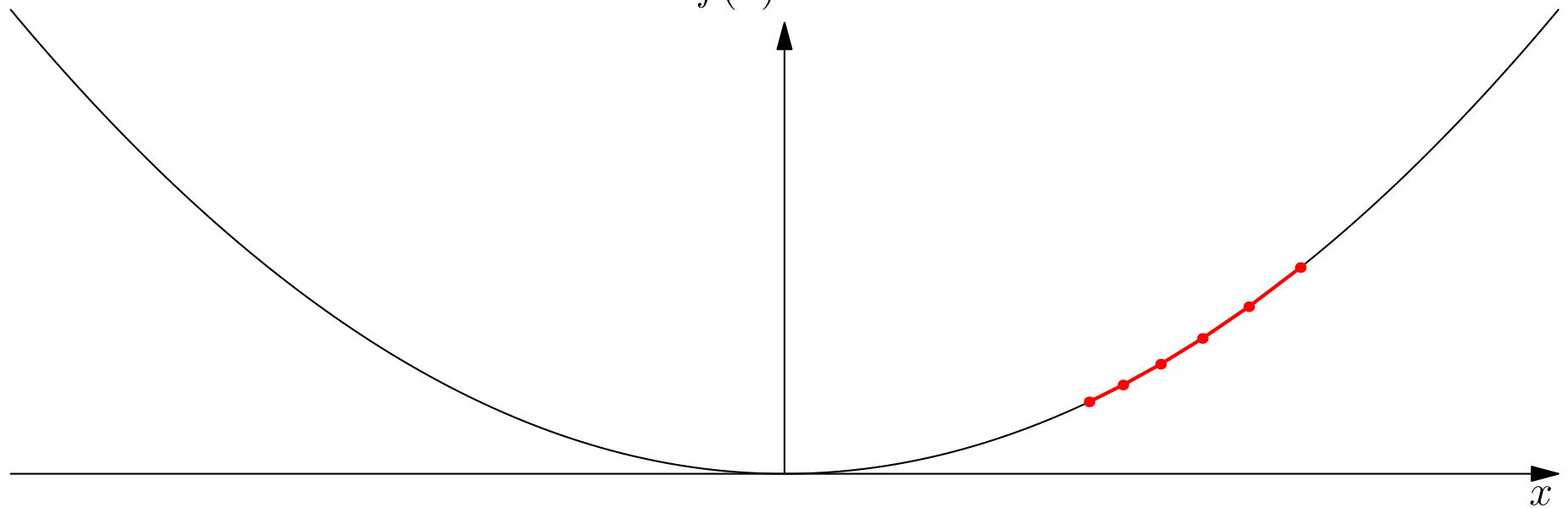


$$r = 0.05$$

# Step Size

$$x \leftarrow x - r f'(x)$$

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

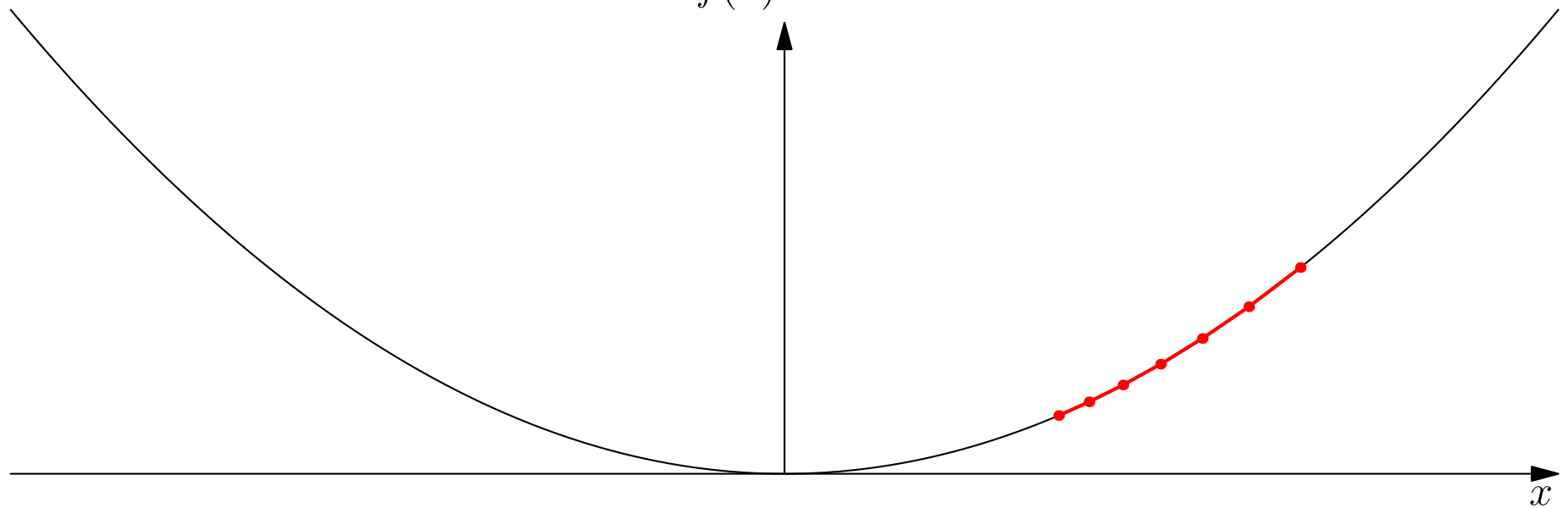


$$r = 0.05$$

# Step Size

$$x \leftarrow x - r f'(x)$$

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

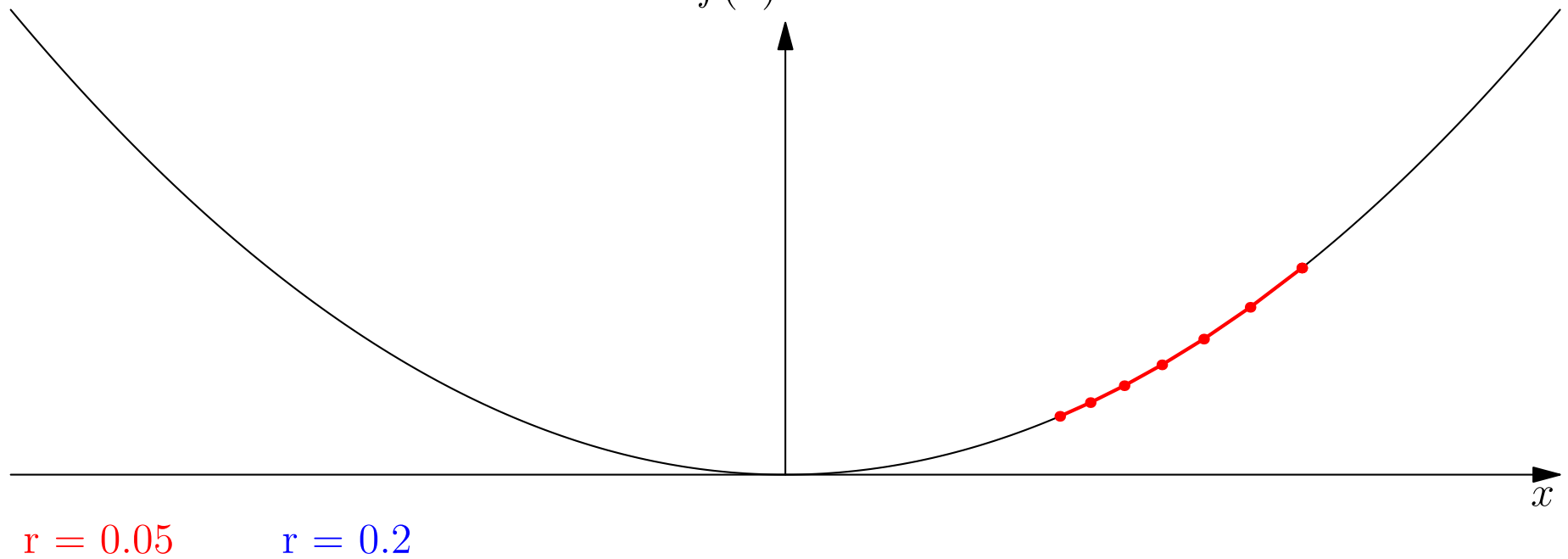


$$r = 0.05$$

# Step Size

$$x \leftarrow x - r f'(x)$$

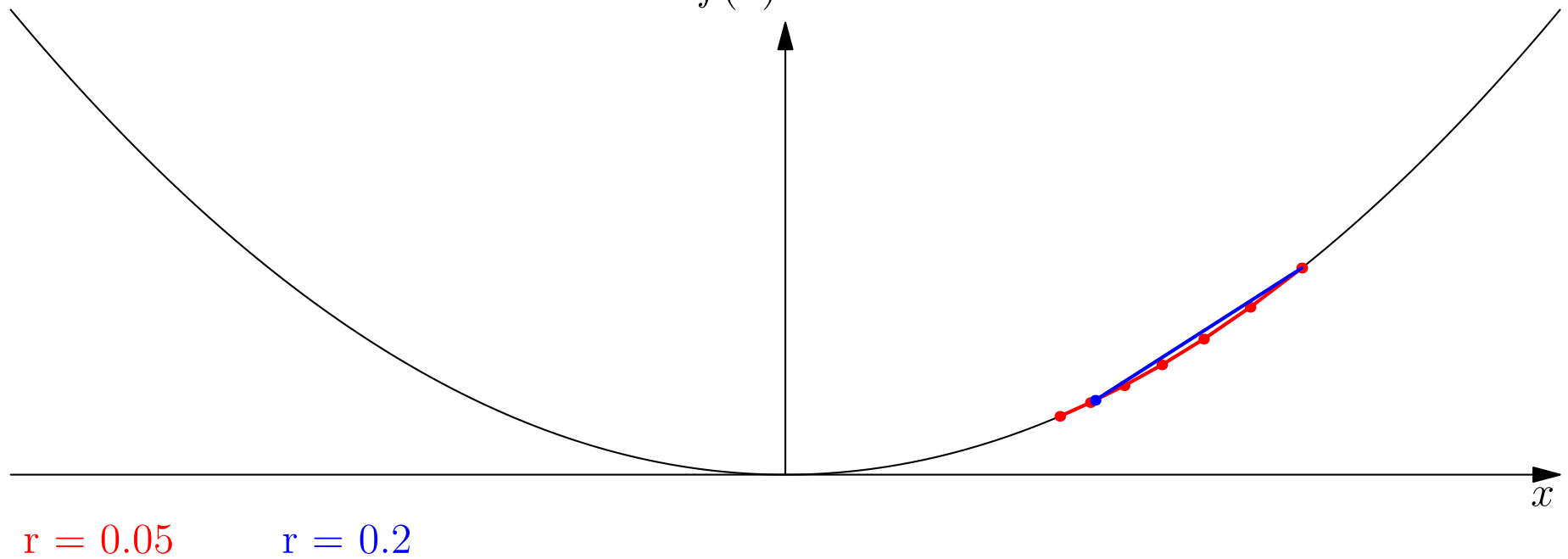
$$f(x) = x^2$$



# Step Size

$$x \leftarrow x - r f'(x)$$

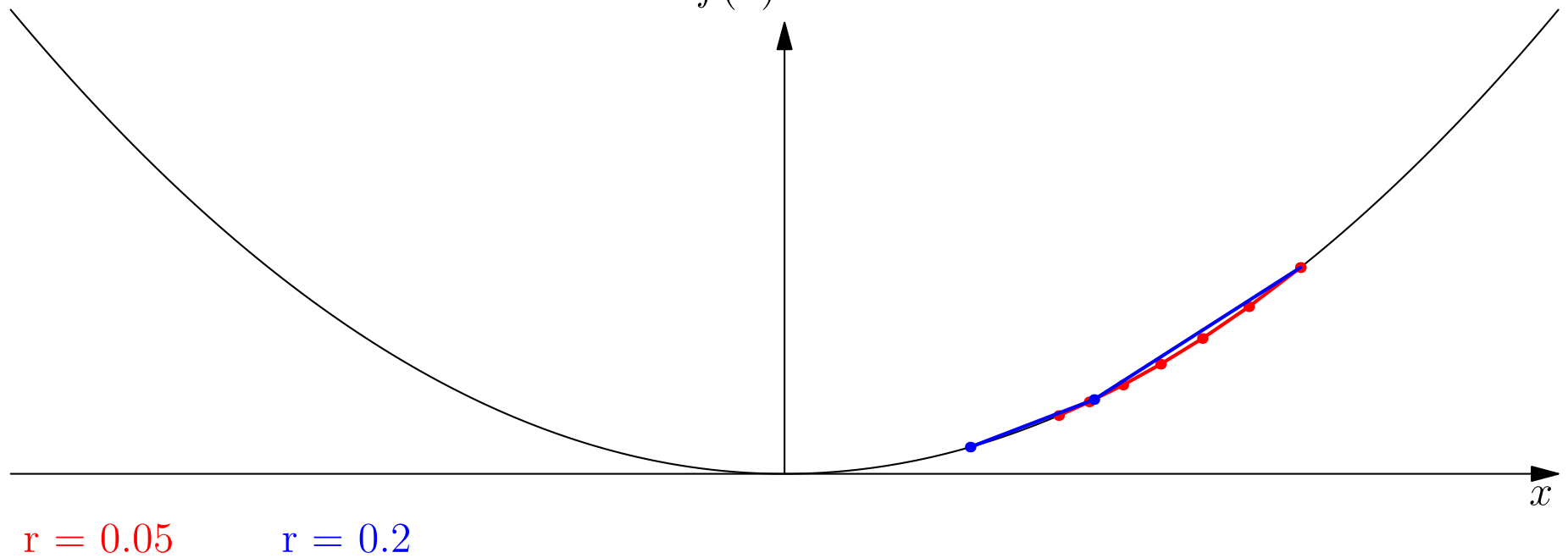
$$f(x) = x^2$$



## Step Size

$$x \leftarrow x - r f'(x)$$

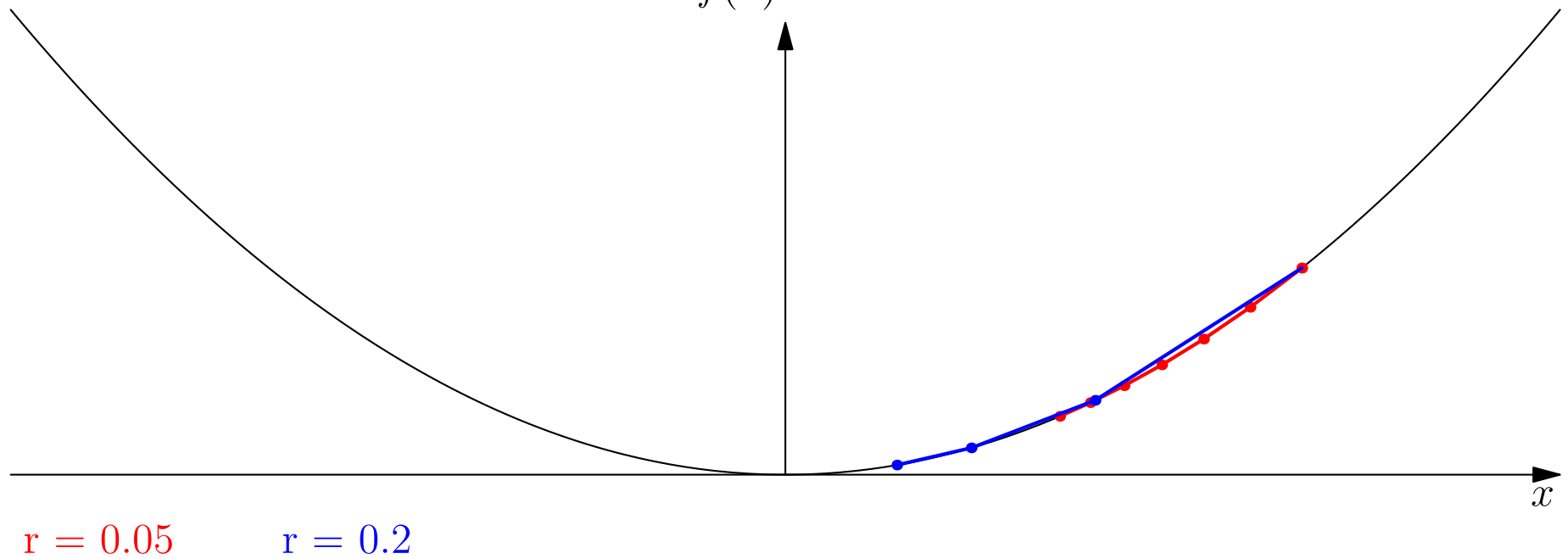
$$f(x) = x^2$$



# Step Size

$$x \leftarrow x - r f'(x)$$

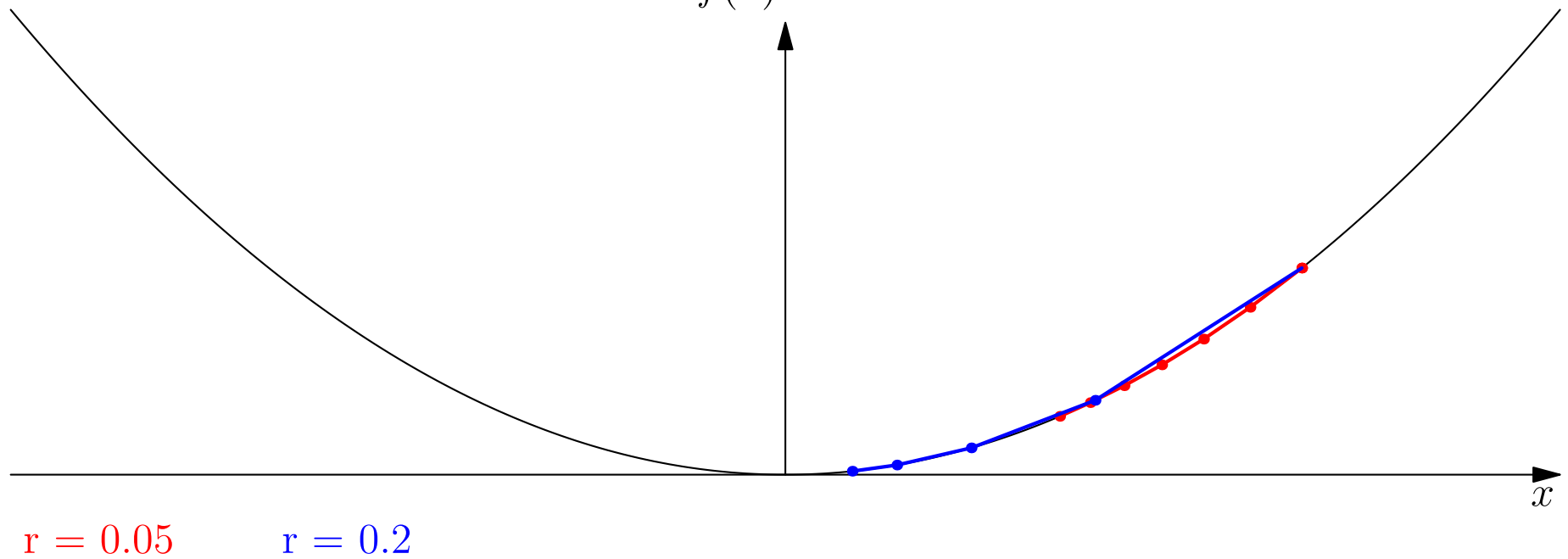
$$f(x) = x^2$$



# Step Size

$$x \leftarrow x - r f'(x)$$

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

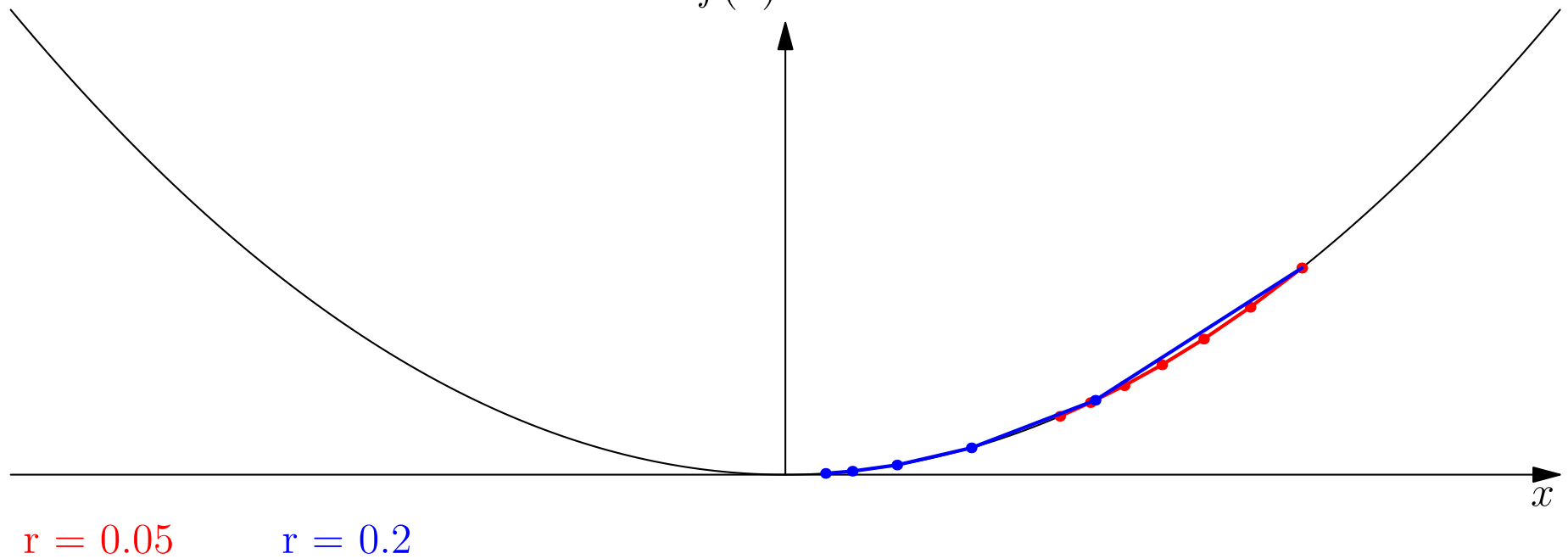




# Step Size

$$x \leftarrow x - r f'(x)$$

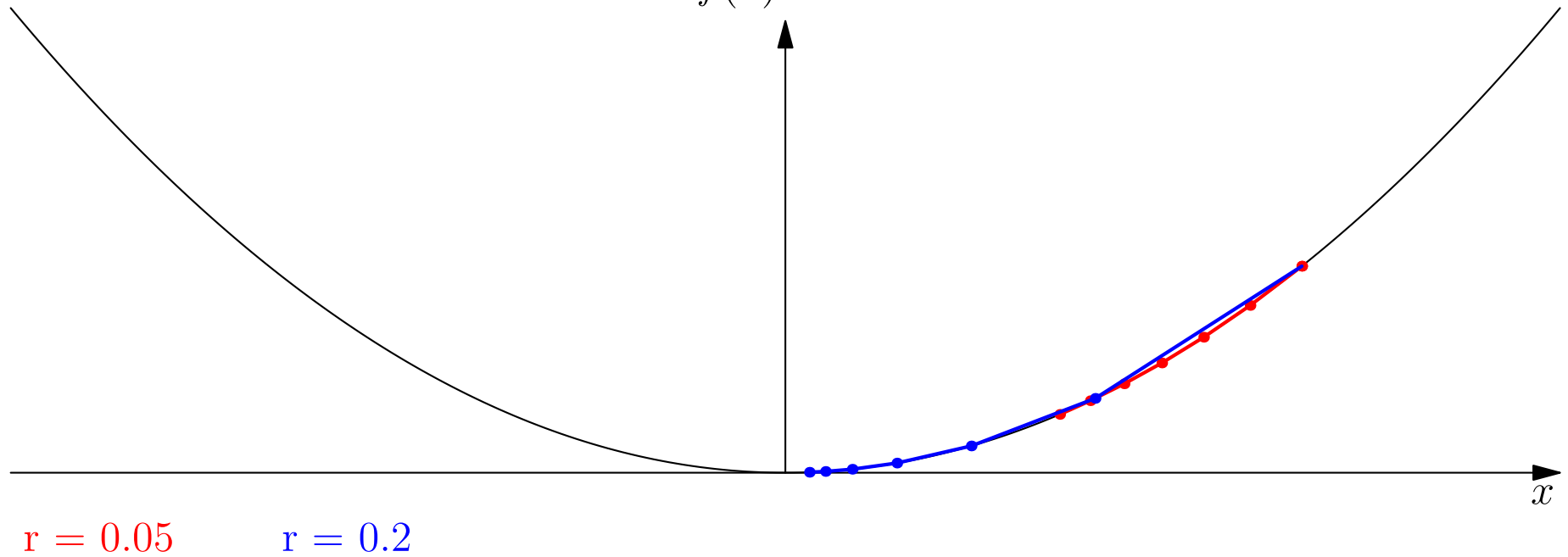
$$f(x) = x^2$$



# Step Size

$$x \leftarrow x - r f'(x)$$

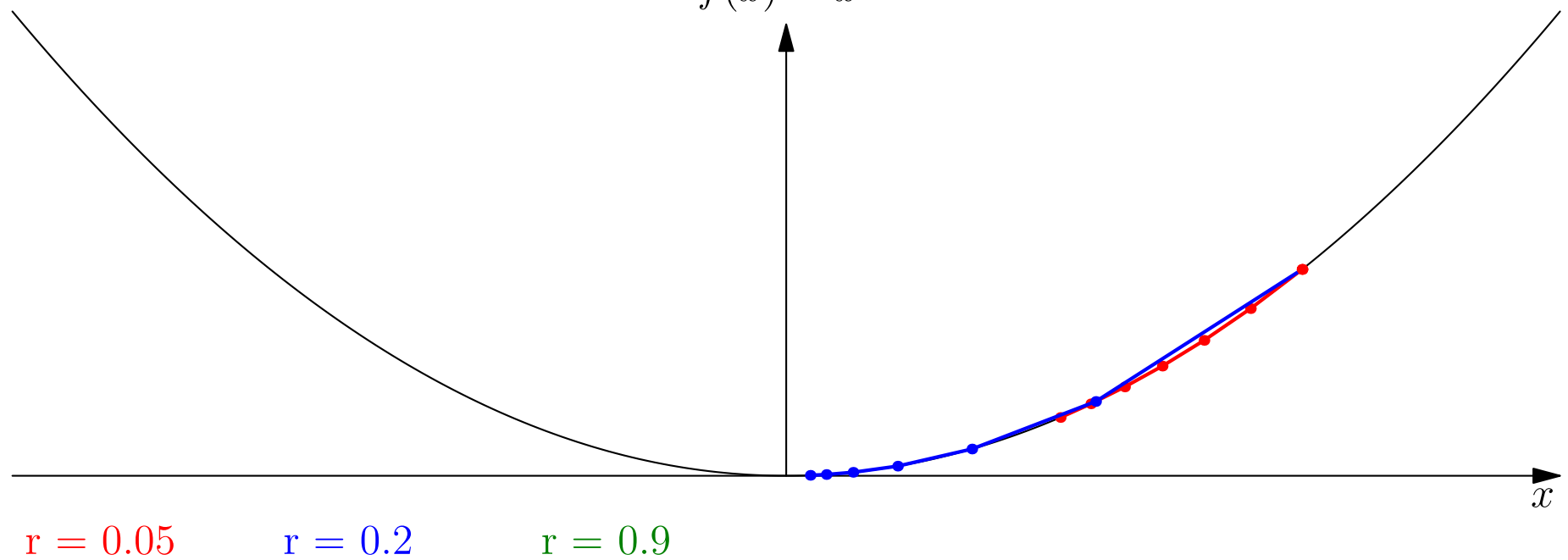
$$f(x) = x^2$$



# Step Size

$$x \leftarrow x - r f'(x)$$

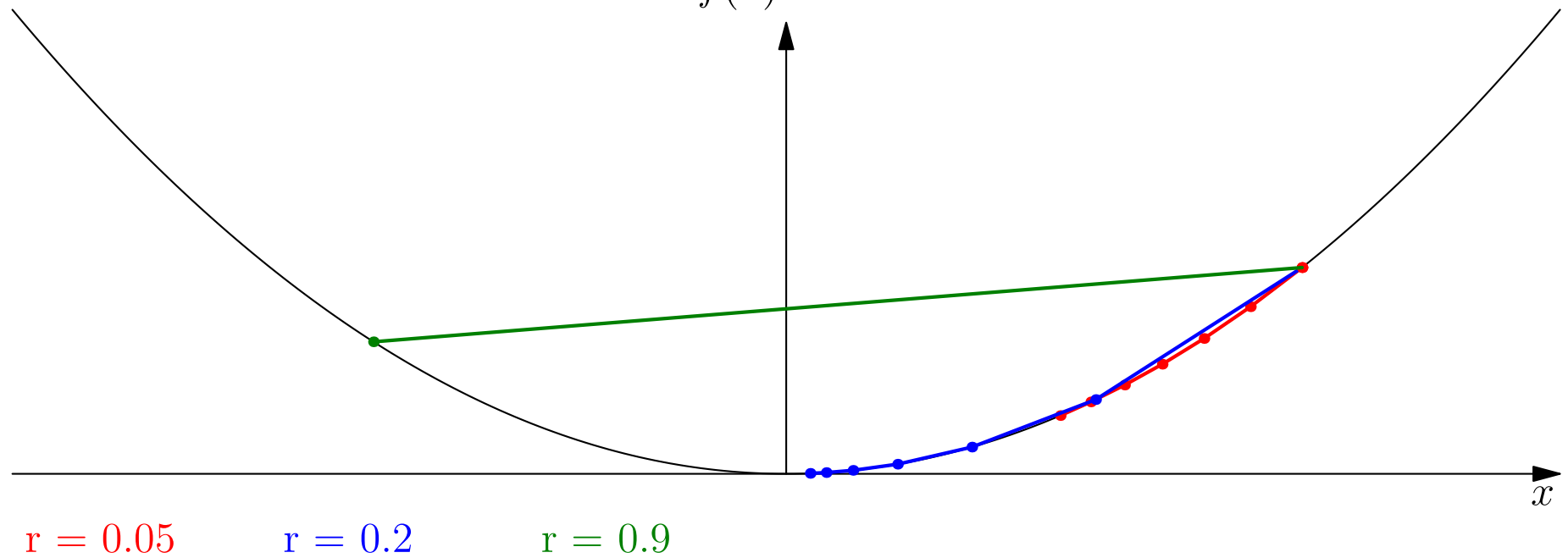
$$f(x) = x^2$$



# Step Size

$$x \leftarrow x - r f'(x)$$

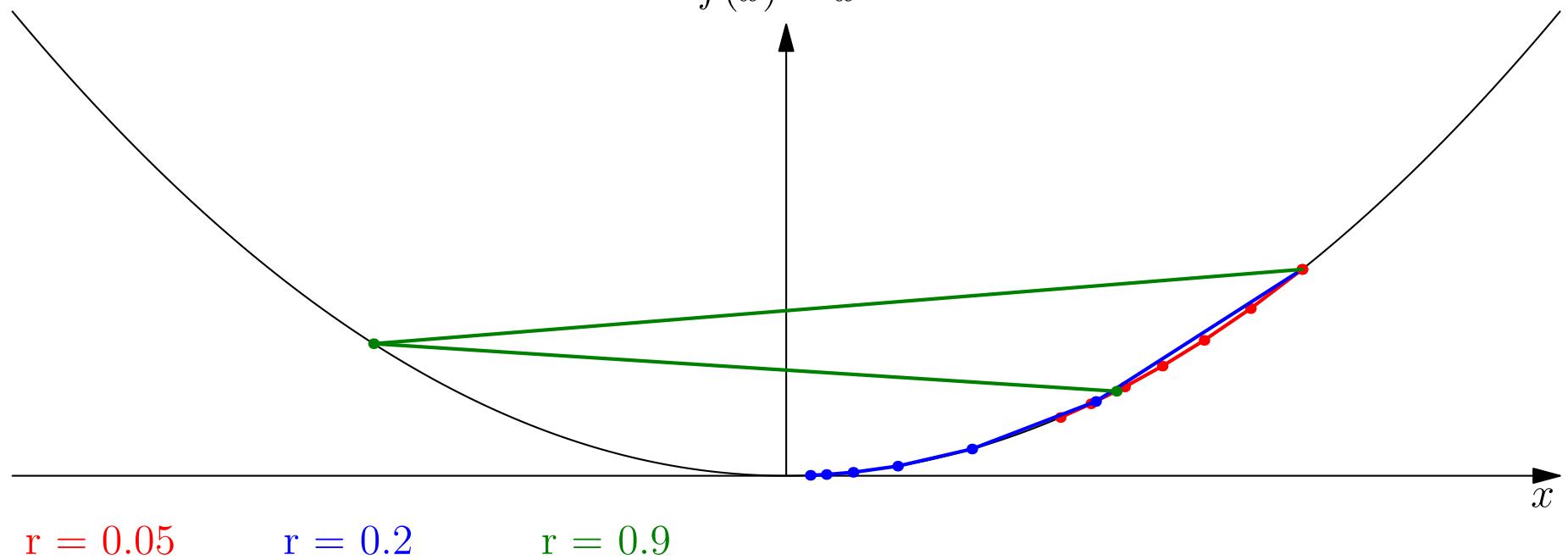
$$f(x) = x^2$$



# Step Size

$$x \leftarrow x - r f'(x)$$

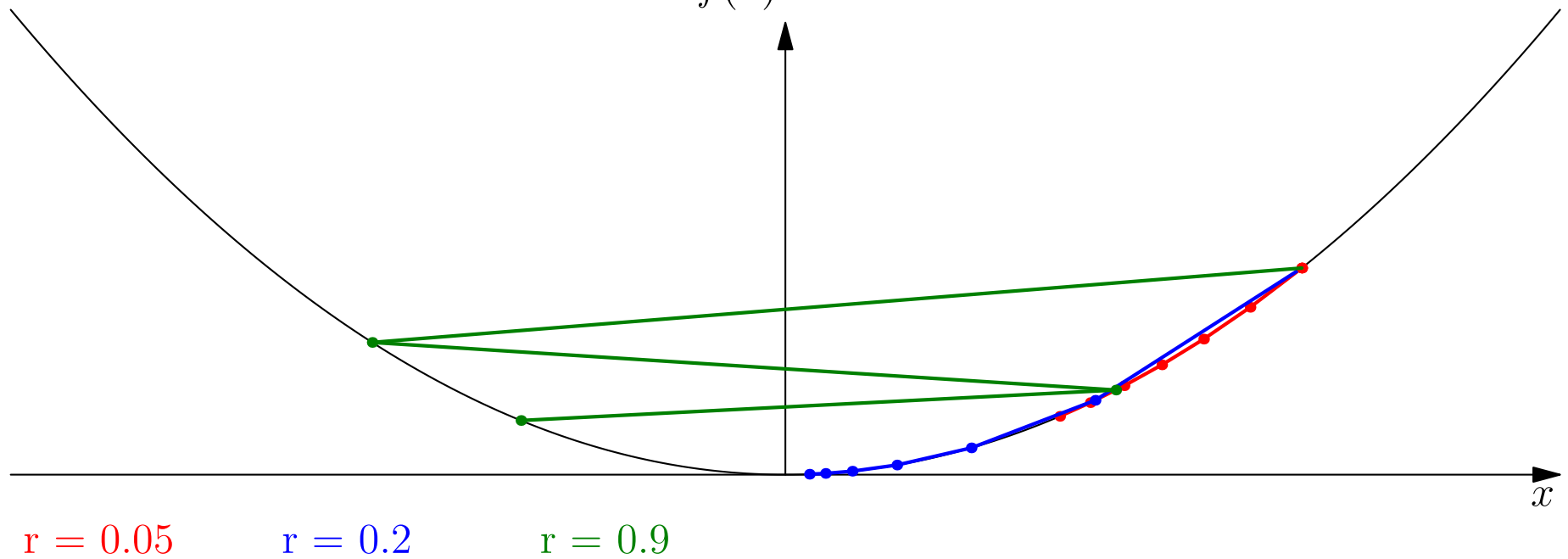
$$f(x) = x^2$$



# Step Size

$$x \leftarrow x - r f'(x)$$

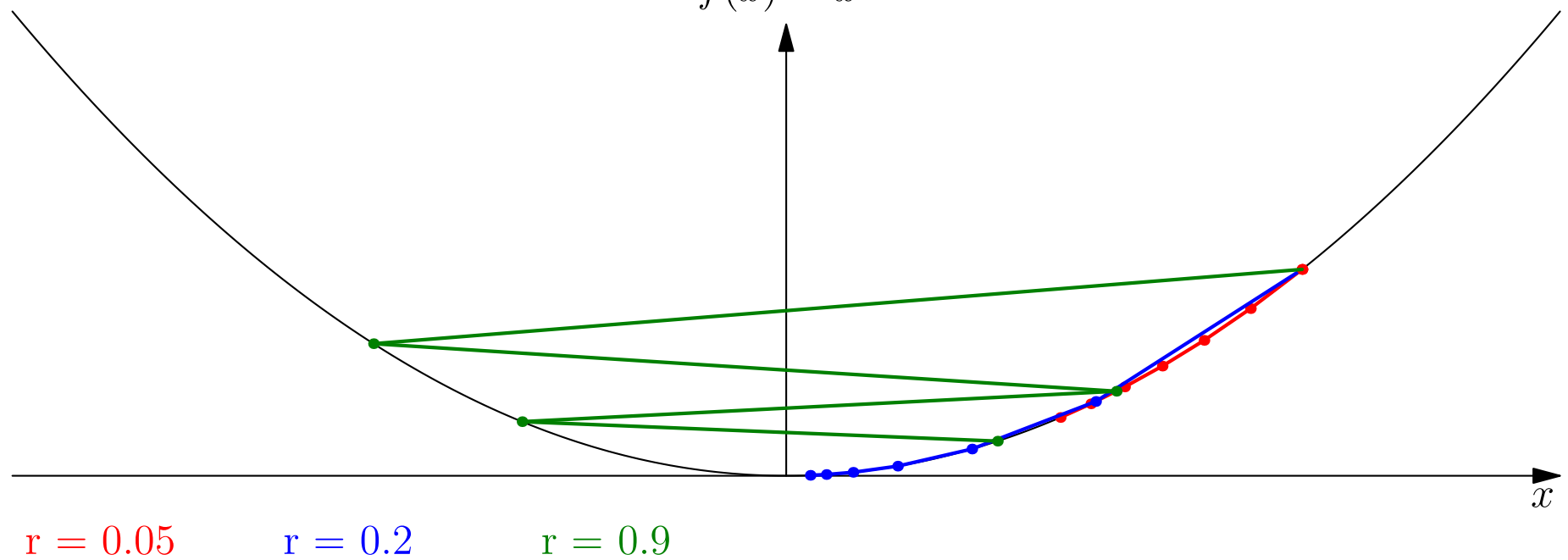
$$f(x) = x^2$$



# Step Size

$$x \leftarrow x - r f'(x)$$

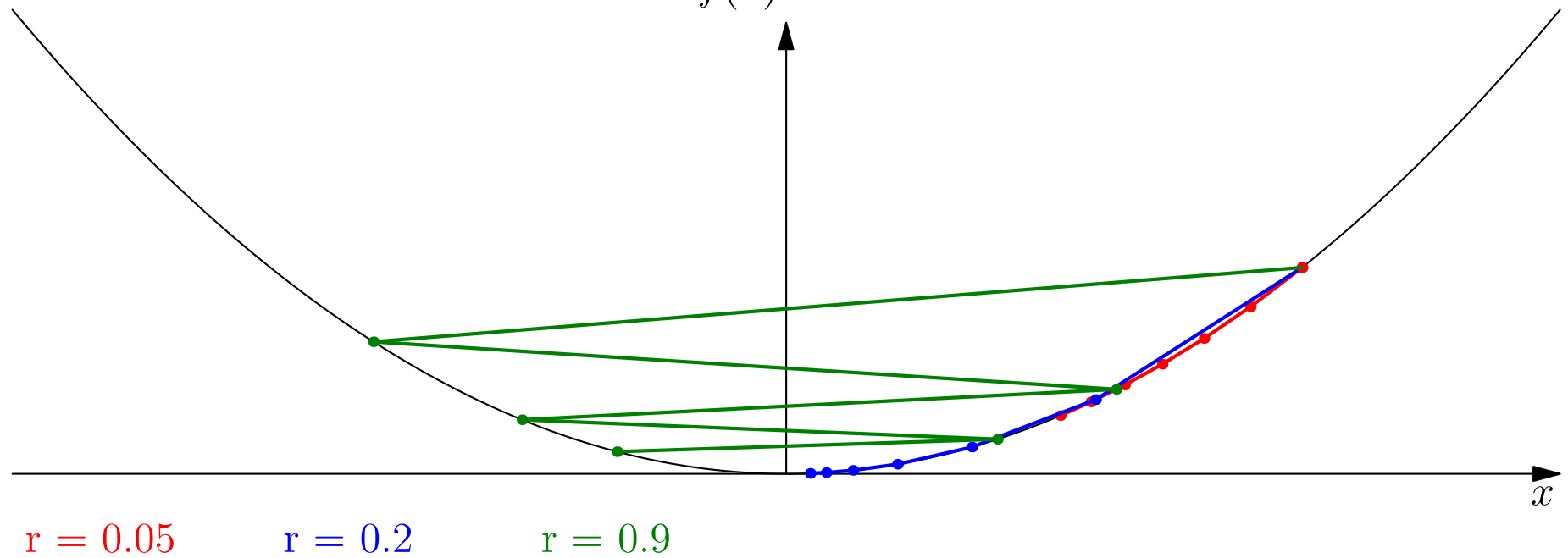
$$f(x) = x^2$$



# Step Size

$$x \leftarrow x - r f'(x)$$

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

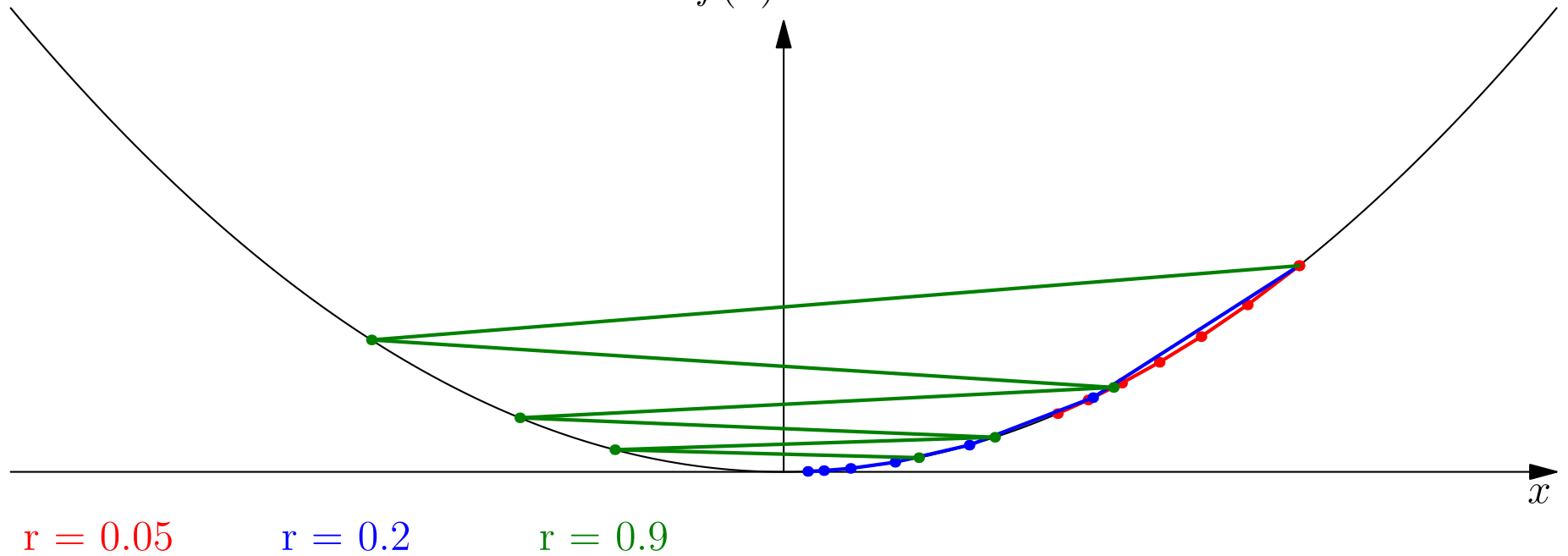




# Step Size

$$x \leftarrow x - r f'(x)$$

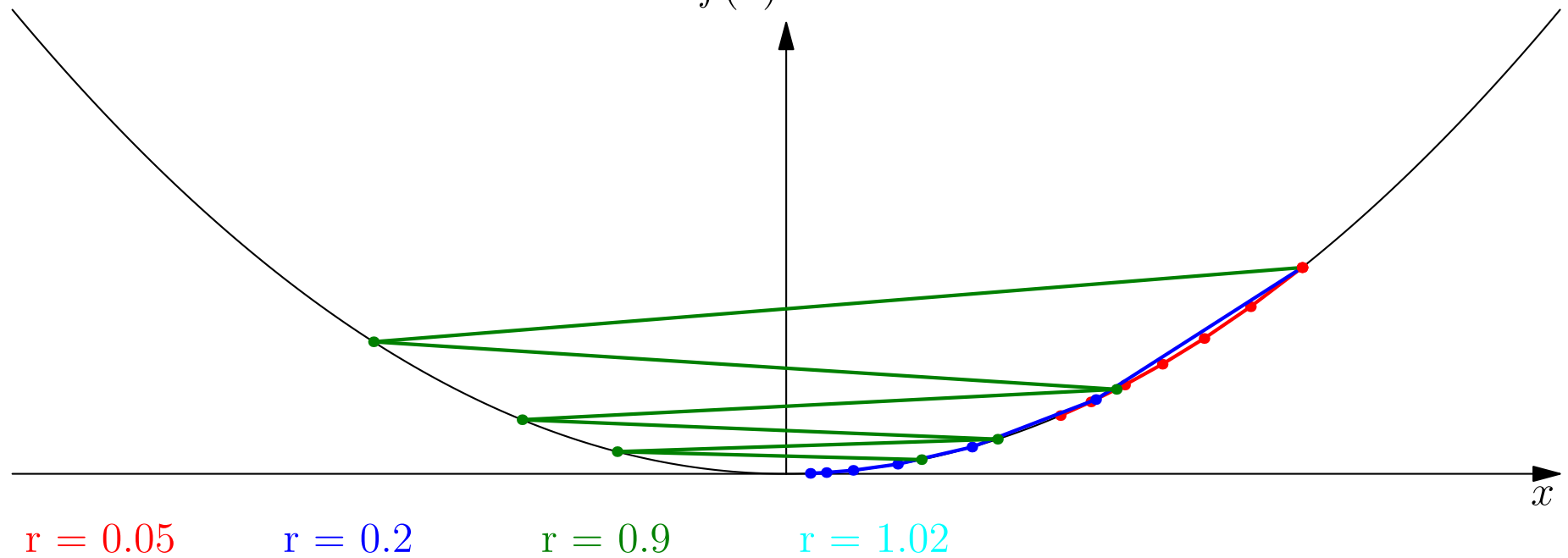
$$f(x) = x^2$$



# Step Size

$$x \leftarrow x - r f'(x)$$

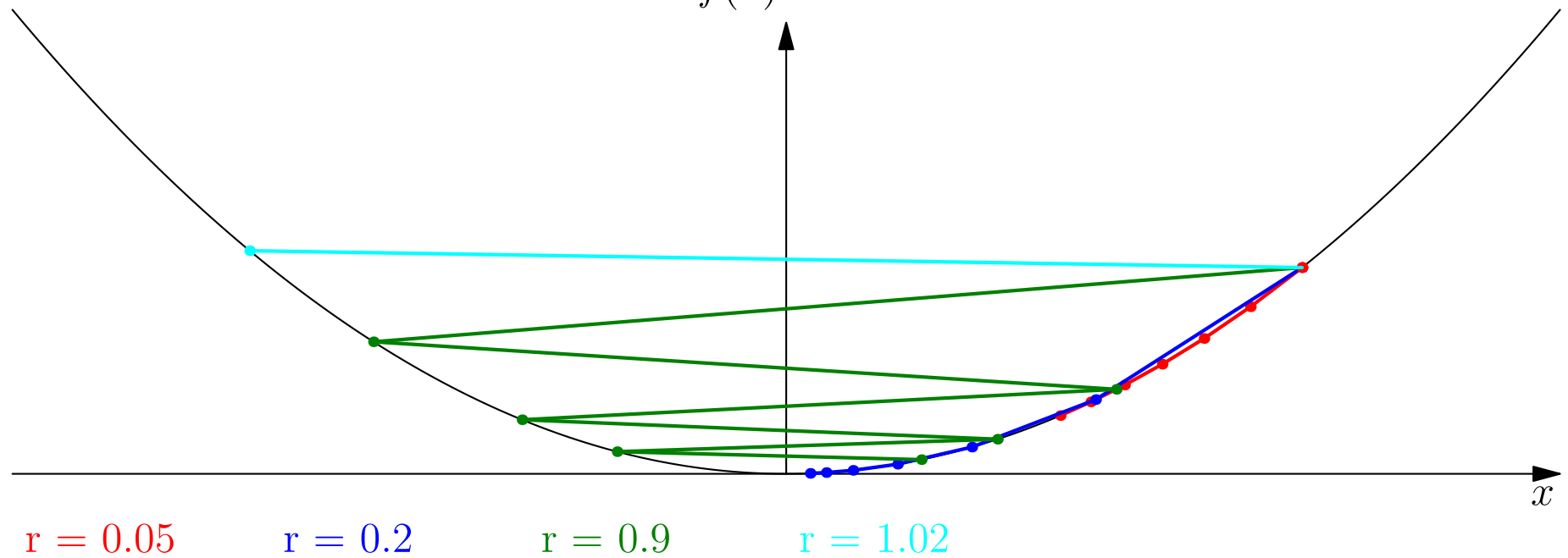
$$f(x) = x^2$$



# Step Size

$$x \leftarrow x - r f'(x)$$

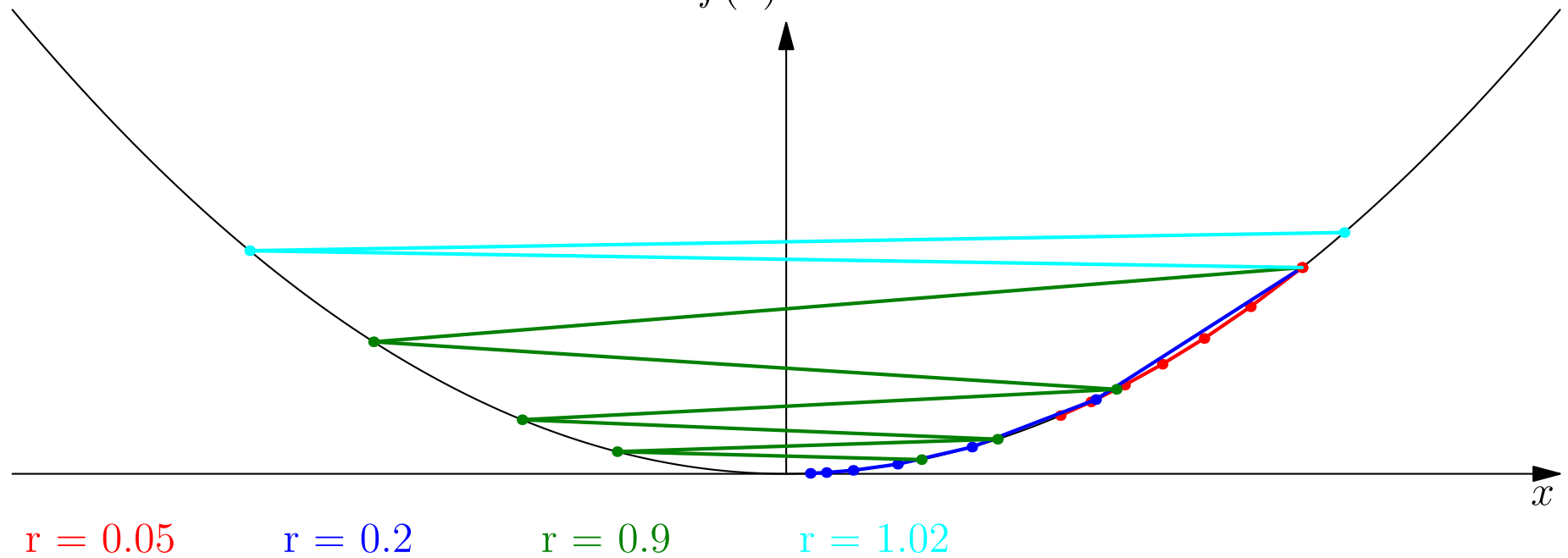
$$f(x) = x^2$$



# Step Size

$$x \leftarrow x - r f'(x)$$

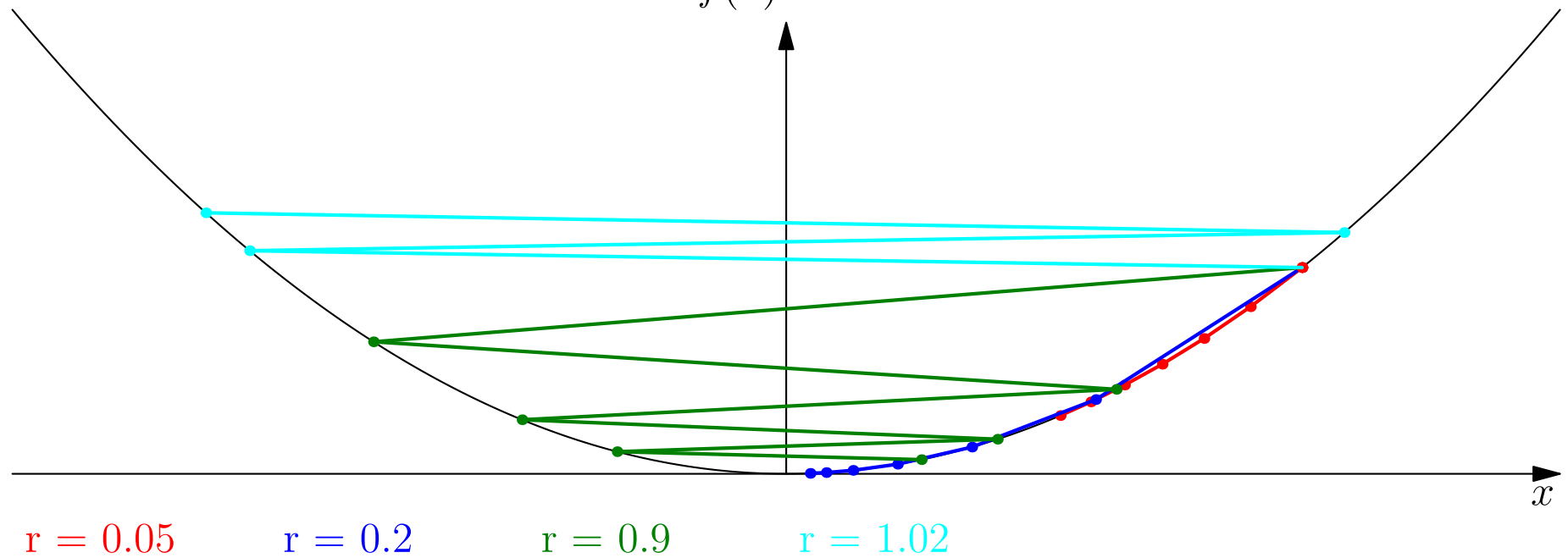
$$f(x) = x^2$$



# Step Size

$$x \leftarrow x - r f'(x)$$

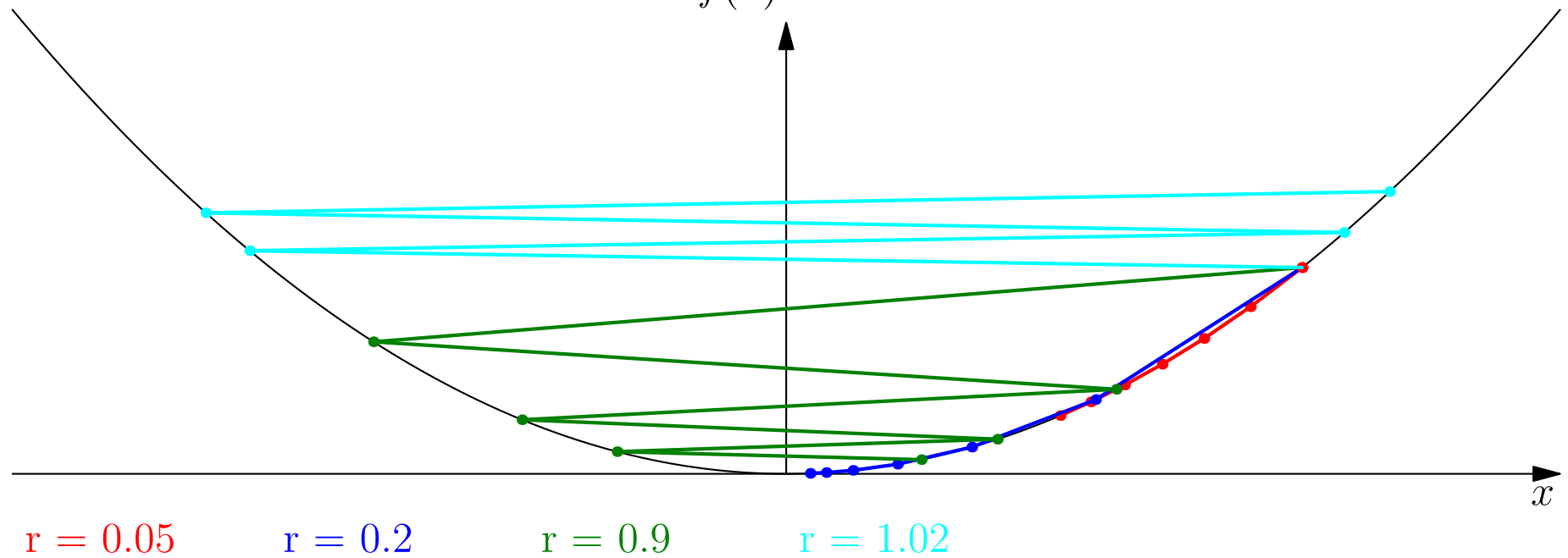
$$f(x) = x^2$$



# Step Size

$$x \leftarrow x - r f'(x)$$

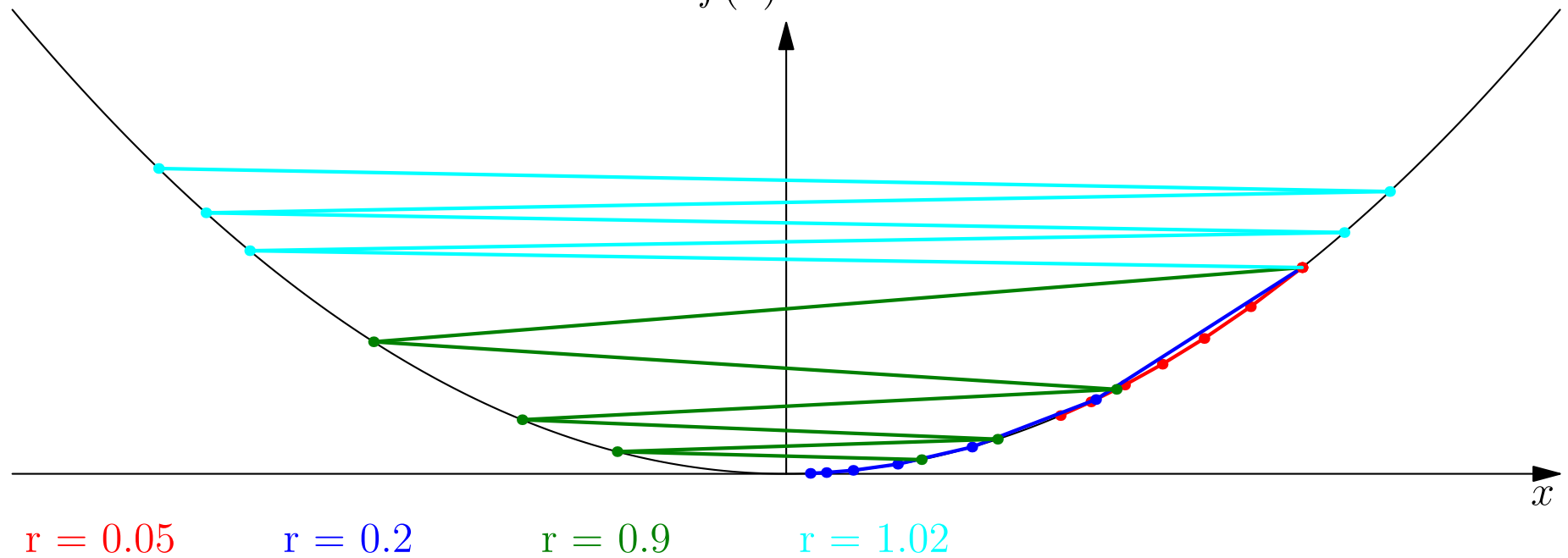
$$f(x) = x^2$$



# Step Size

$$x \leftarrow x - r f'(x)$$

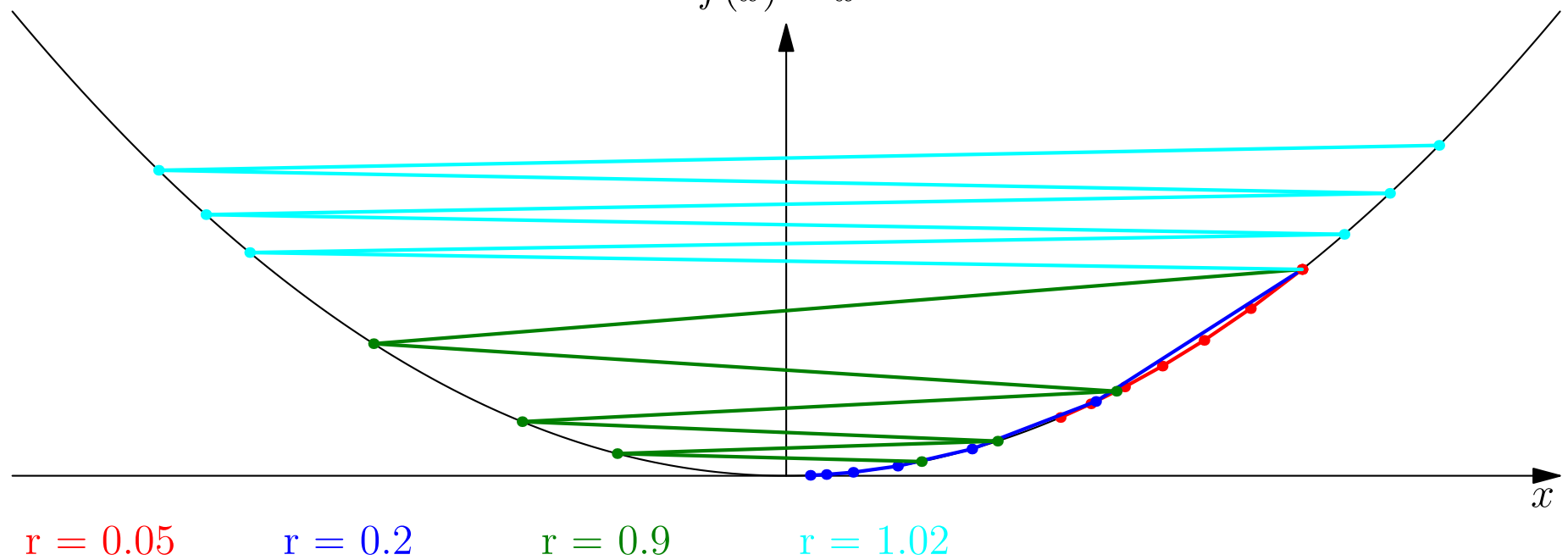
$$f(x) = x^2$$



# Step Size

$$x \leftarrow x - r f'(x)$$

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



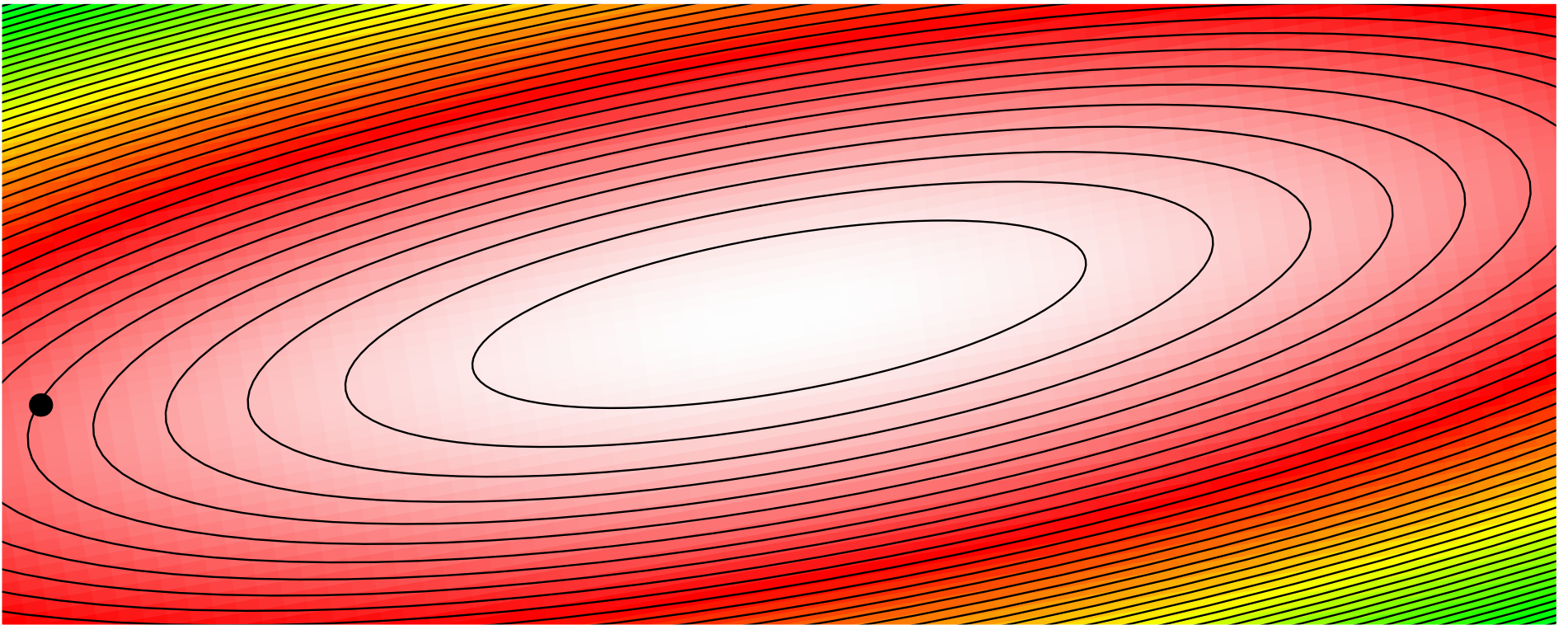


# Higher Dimensions

- In higher dimensions the problem is that there are some directions you need to move a long way

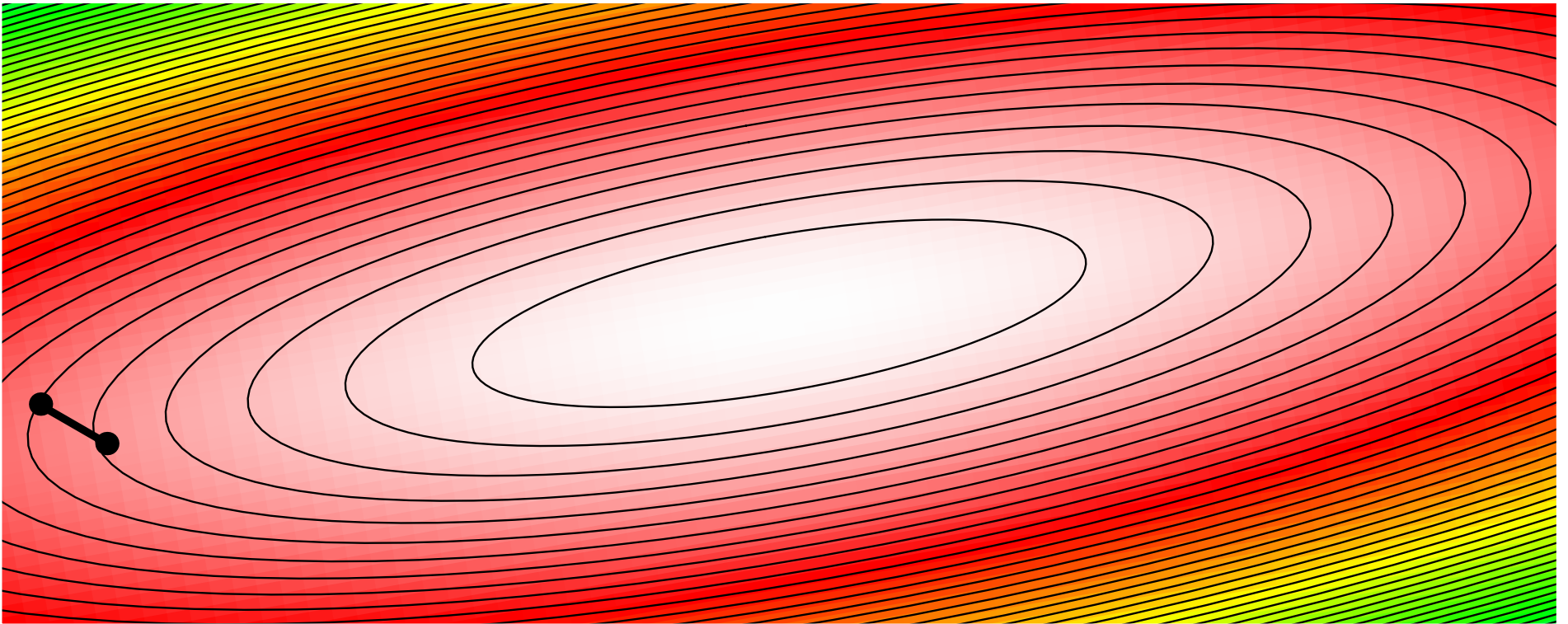
# Higher Dimensions

- In higher dimensions the problem is that there are some directions you need to move a long way



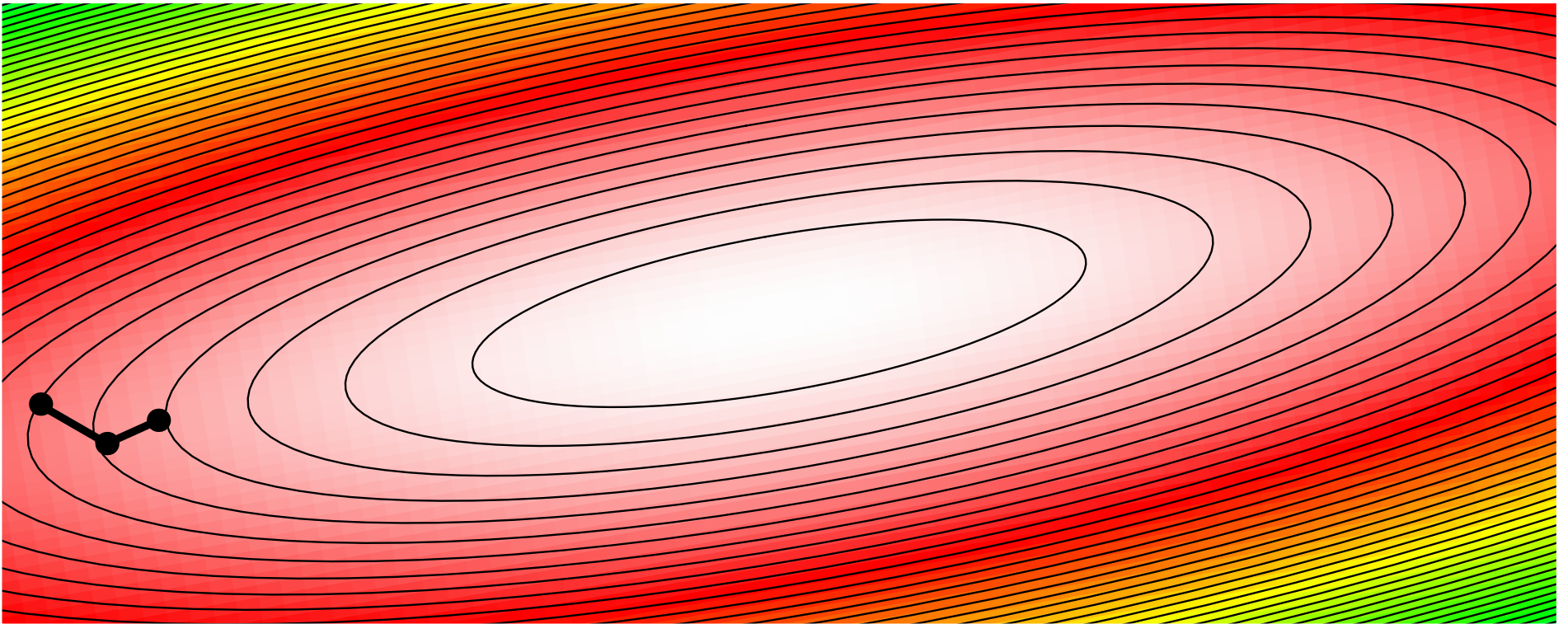
# Higher Dimensions

- In higher dimensions the problem is that there are some directions you need to move a long way



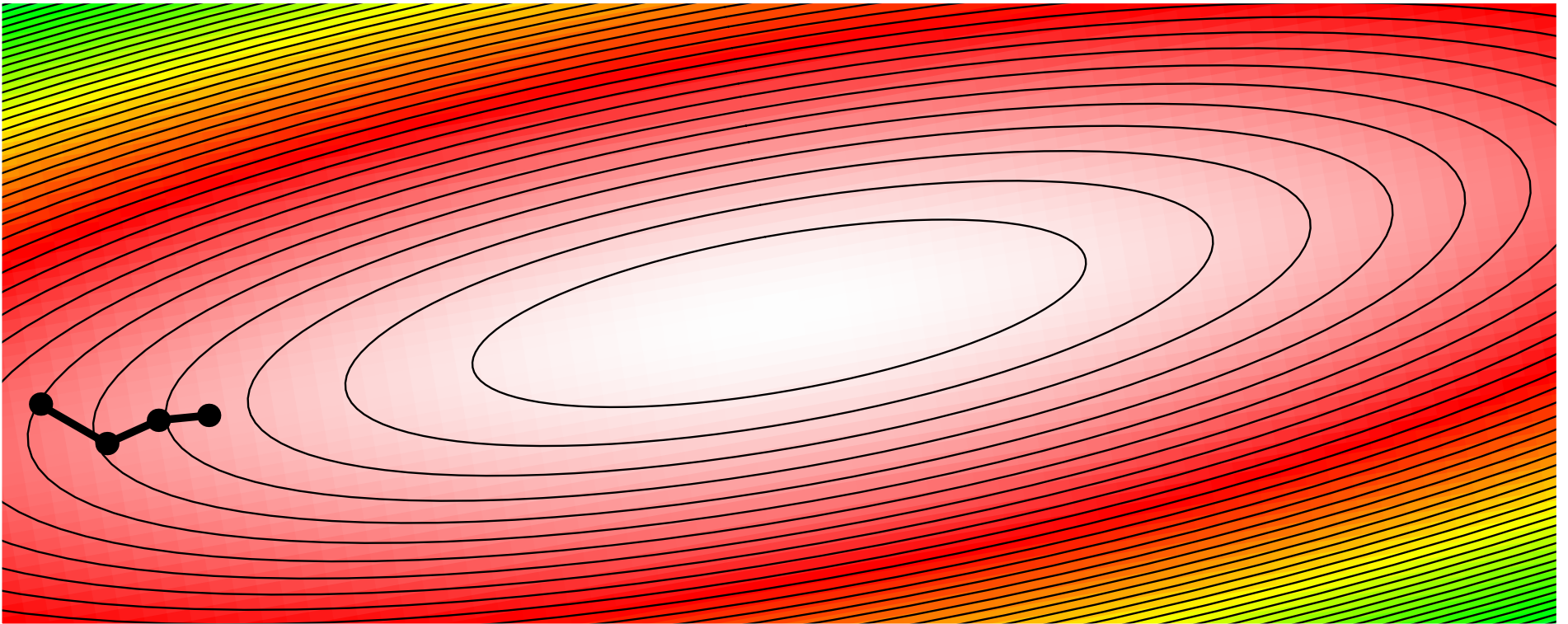
# Higher Dimensions

- In higher dimensions the problem is that there are some directions you need to move a long way



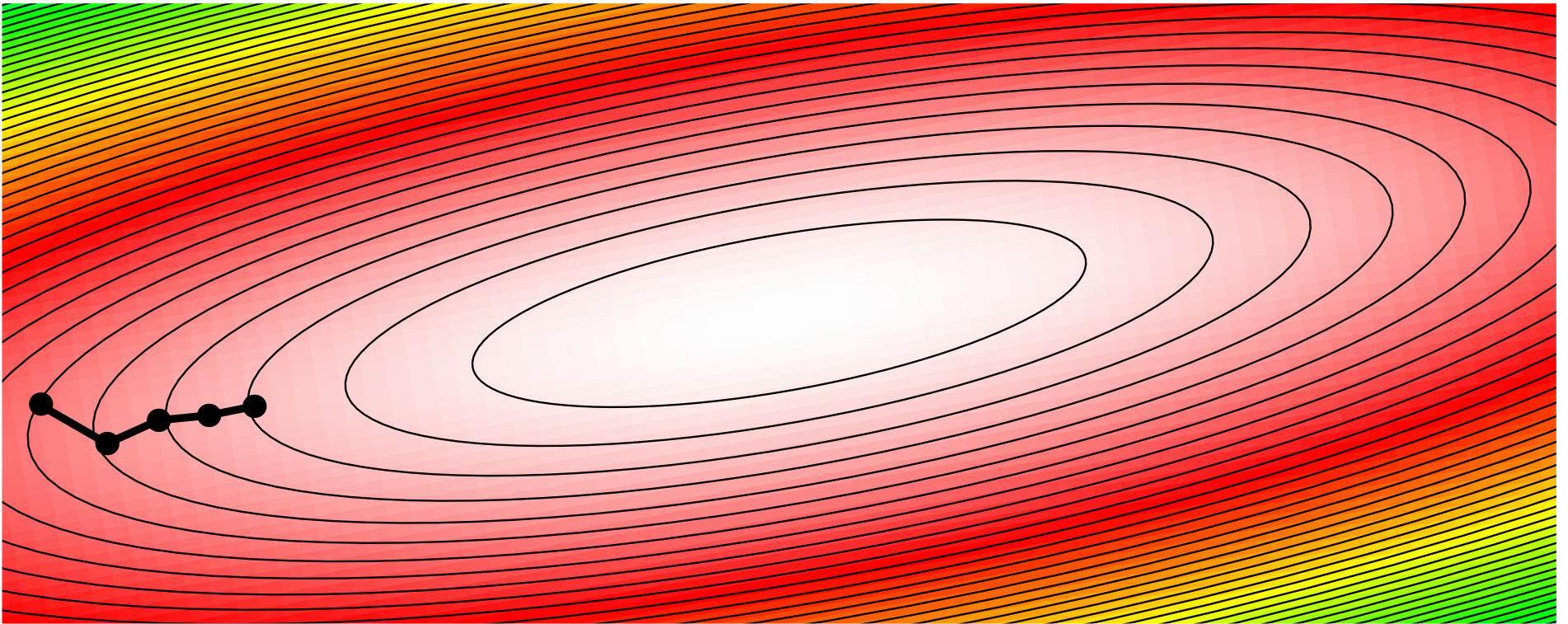
# Higher Dimensions

- In higher dimensions the problem is that there are some directions you need to move a long way



# Higher Dimensions

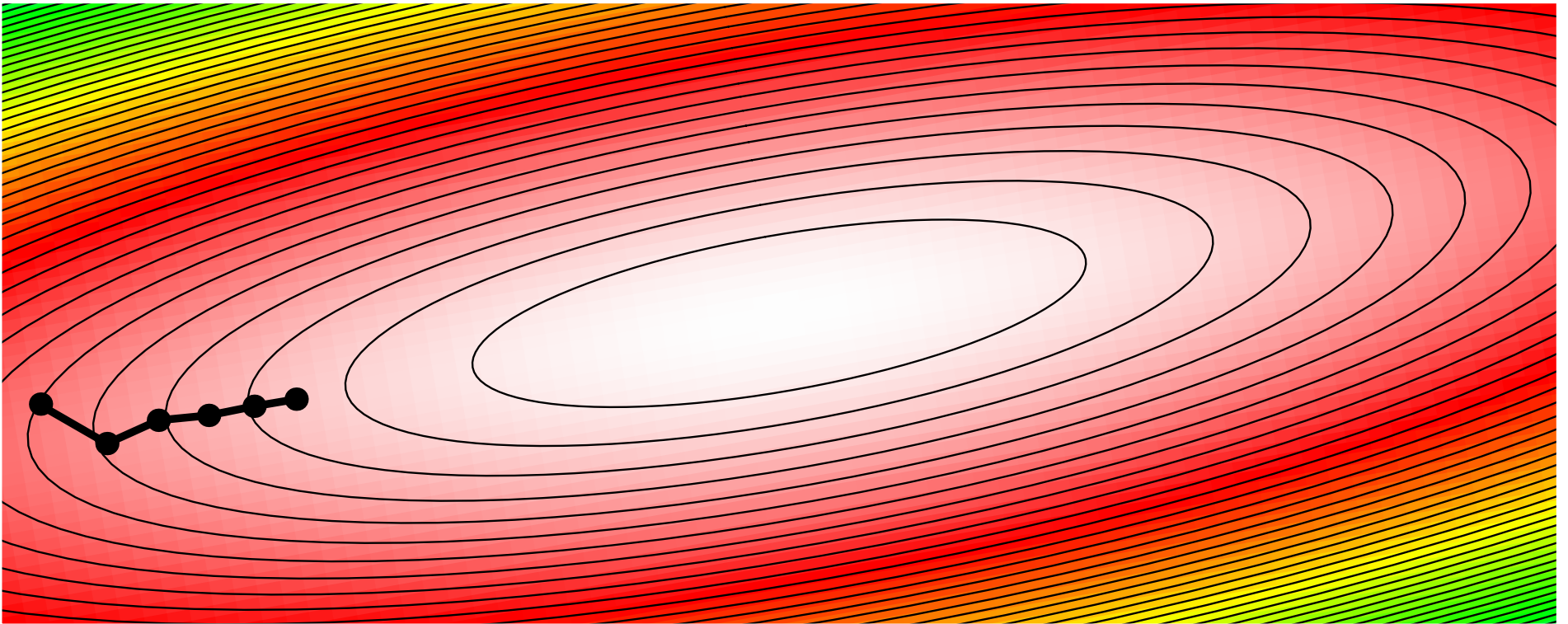
- In higher dimensions the problem is that there are some directions you need to move a long way





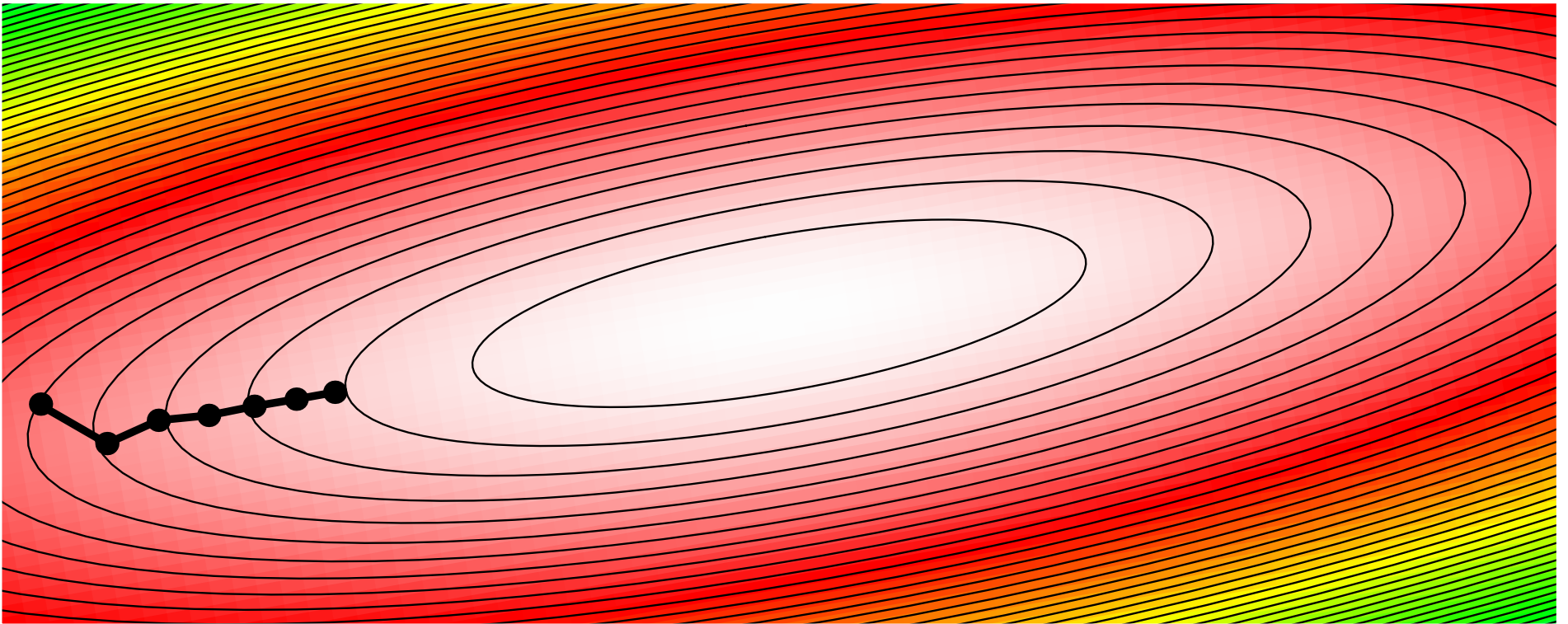
# Higher Dimensions

- In higher dimensions the problem is that there are some directions you need to move a long way



# Higher Dimensions

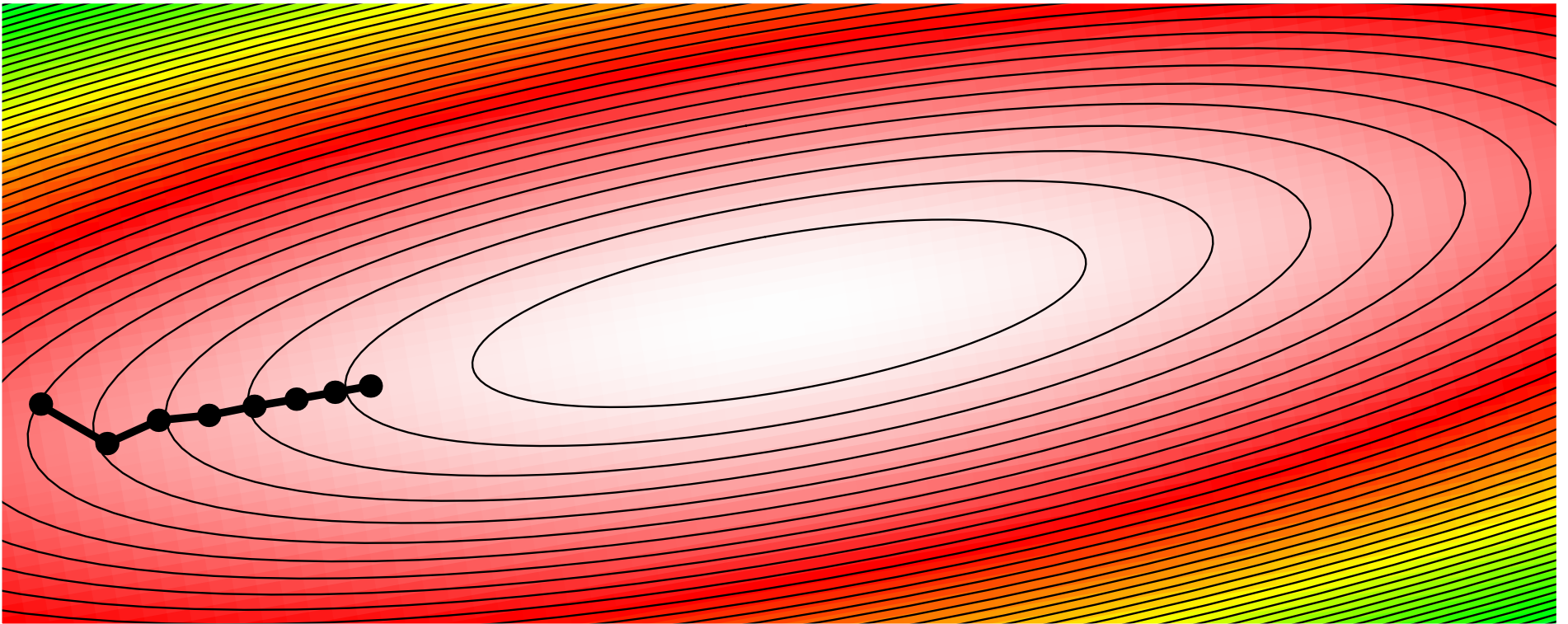
- In higher dimensions the problem is that there are some directions you need to move a long way





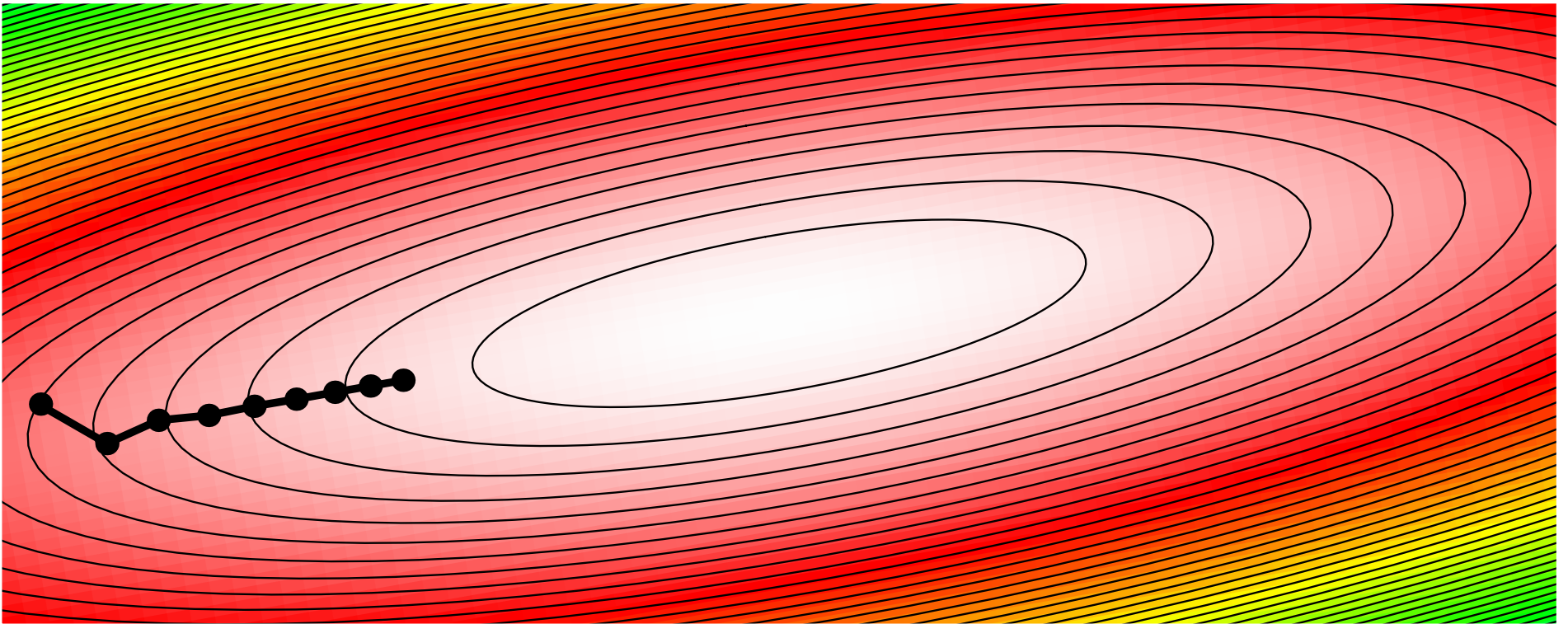
# Higher Dimensions

- In higher dimensions the problem is that there are some directions you need to move a long way



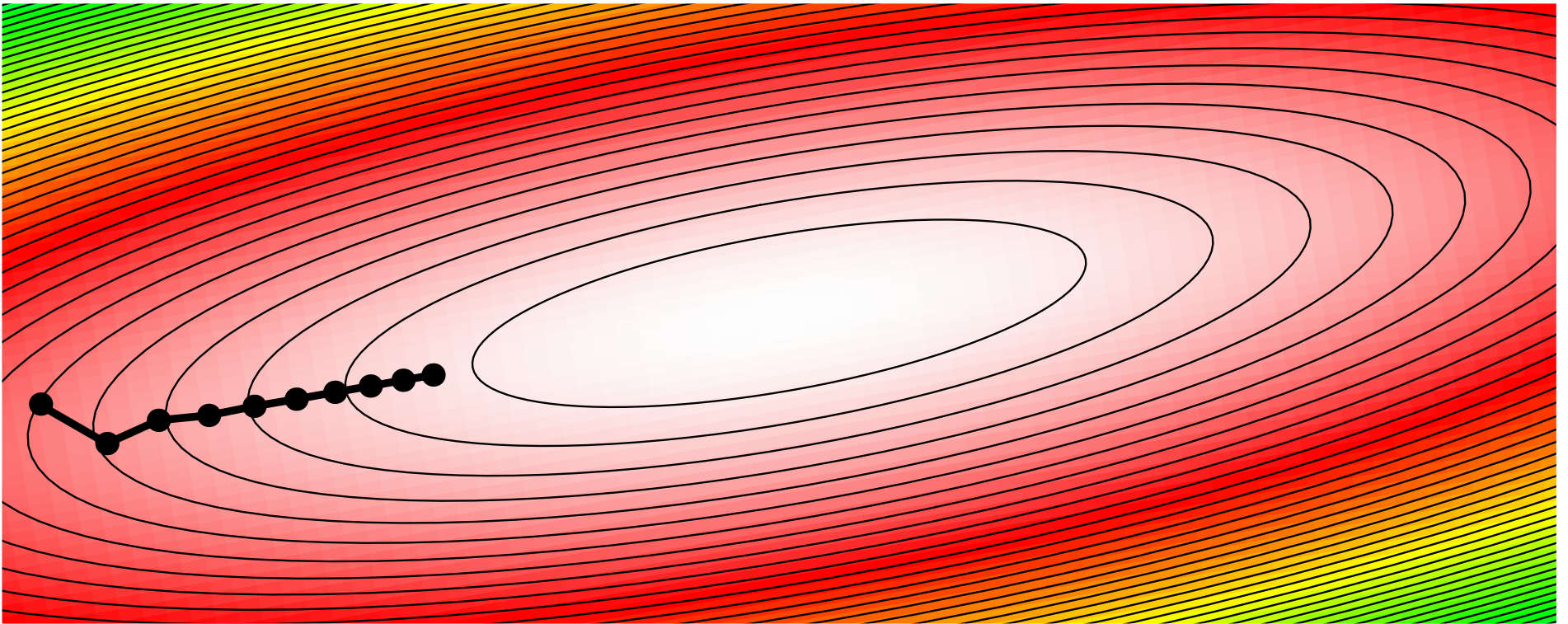
# Higher Dimensions

- In higher dimensions the problem is that there are some directions you need to move a long way



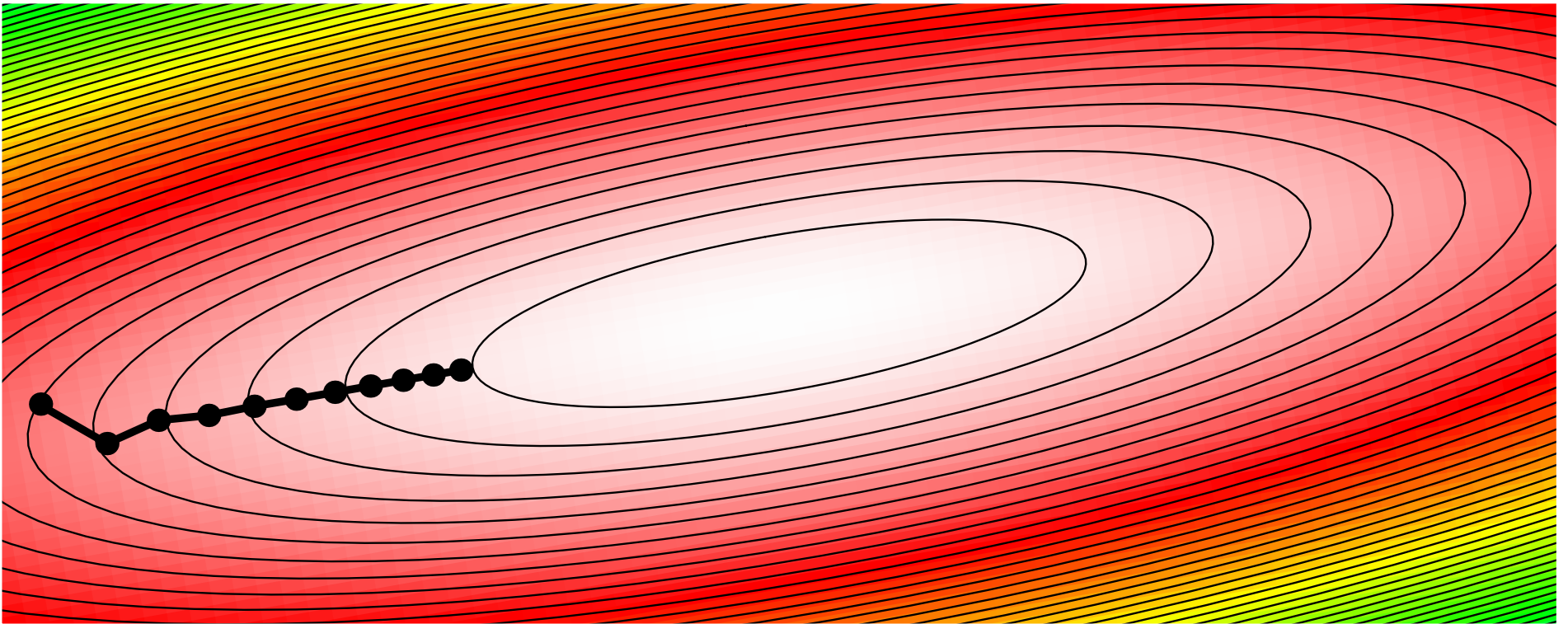
# Higher Dimensions

- In higher dimensions the problem is that there are some directions you need to move a long way



# Higher Dimensions

- In higher dimensions the problem is that there are some directions you need to move a long way

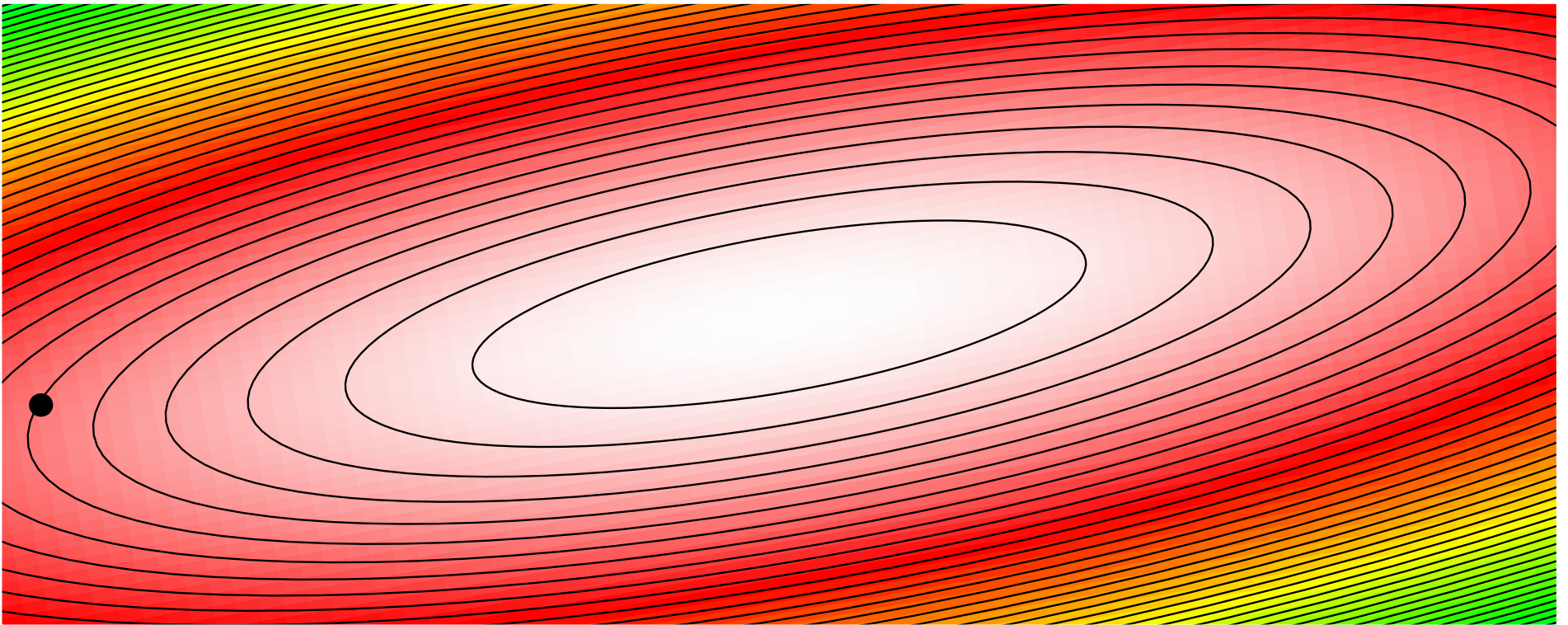


# Getting There Quicker

- Increasing the step size speeds up convergence, but the direction of steepest descent doesn't point to the minimum

# Getting There Quicker

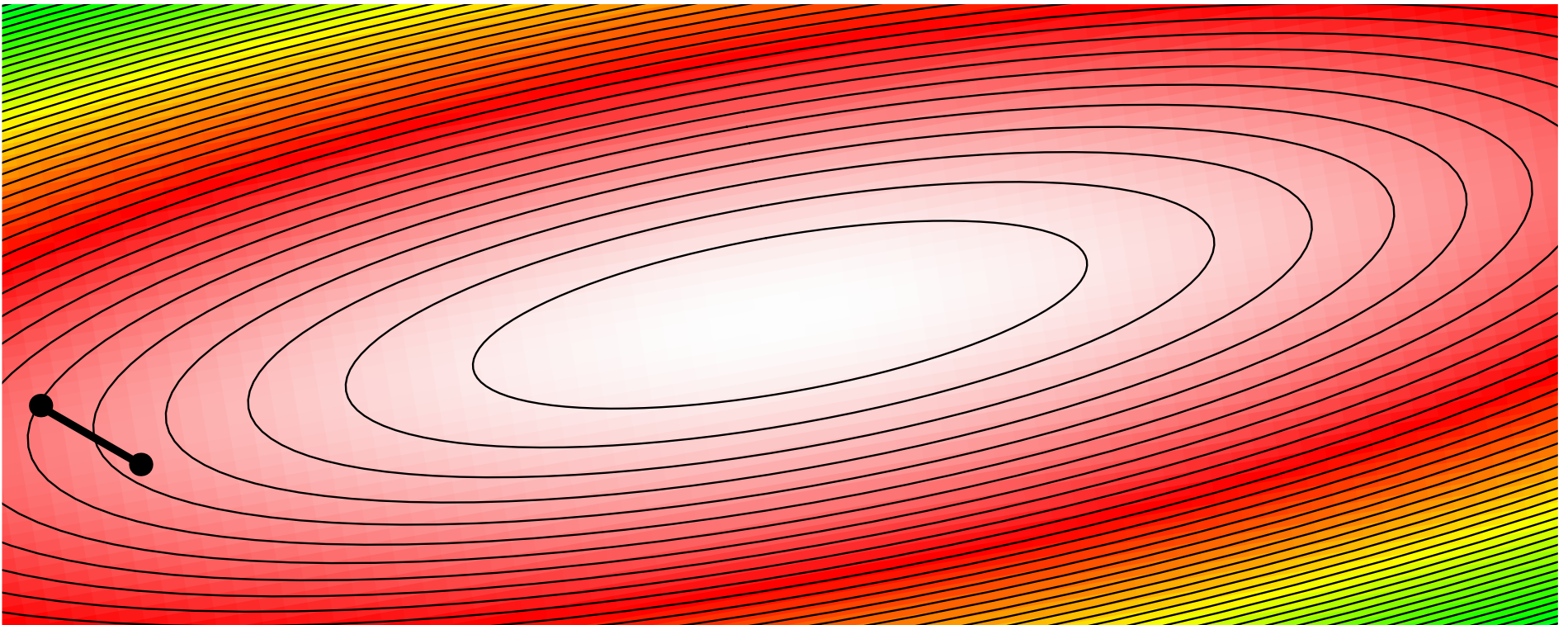
- Increasing the step size speeds up convergence, but the direction of steepest descent doesn't point to the minimum





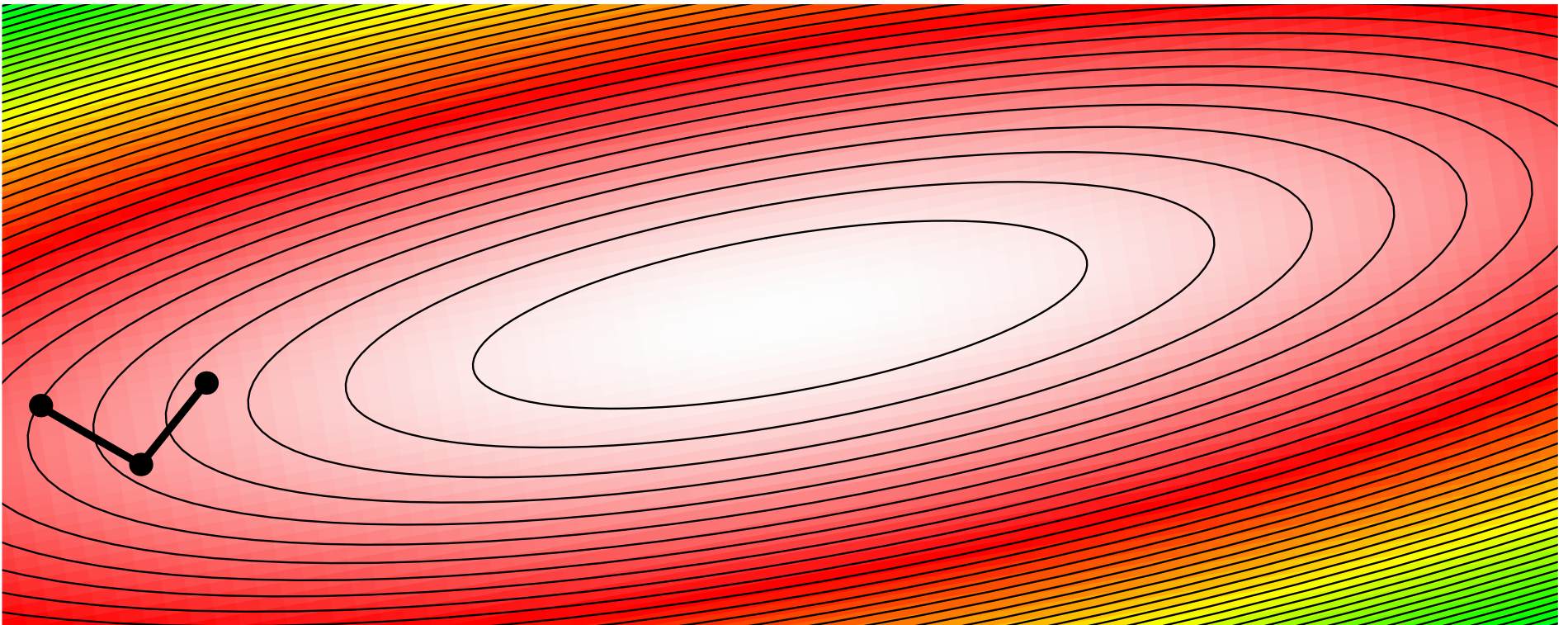
# Getting There Quicker

- Increasing the step size speeds up convergence, but the direction of steepest descent doesn't point to the minimum



# Getting There Quicker

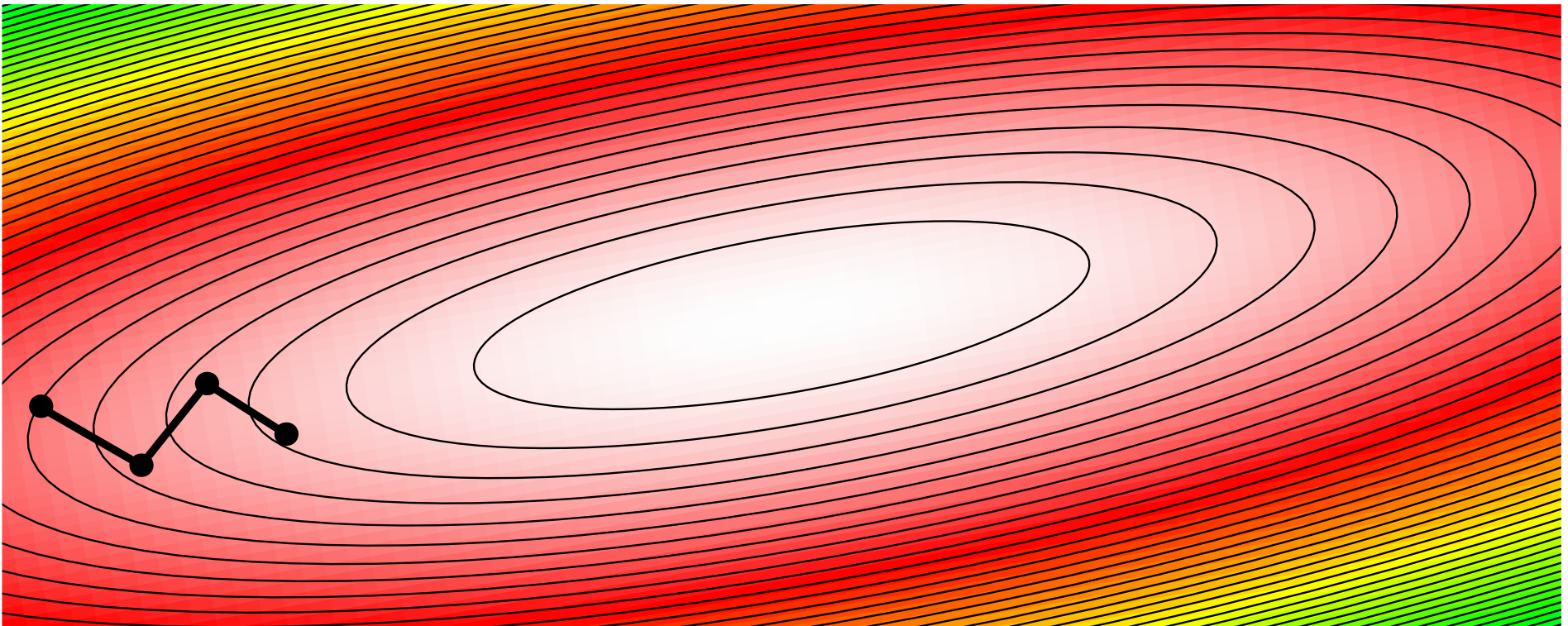
- Increasing the step size speeds up convergence, but the direction of steepest descent doesn't point to the minimum





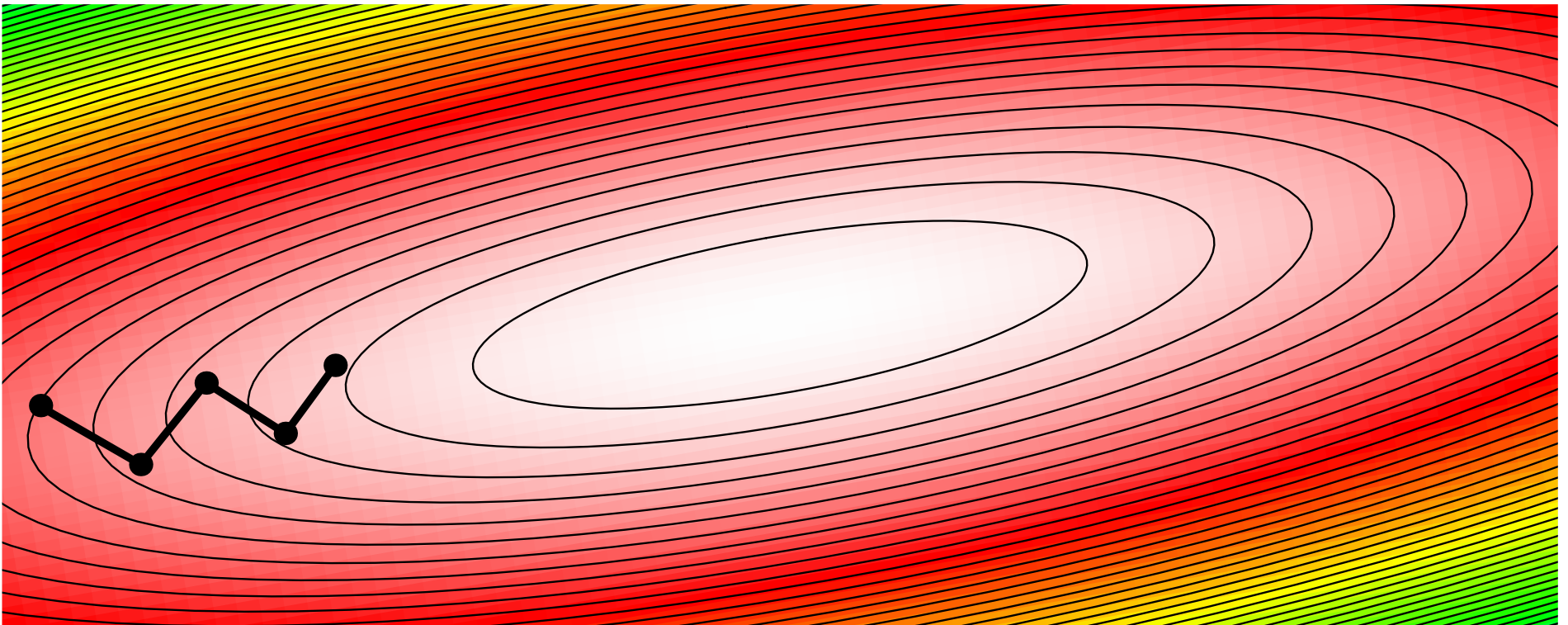
# Getting There Quicker

- Increasing the step size speeds up convergence, but the direction of steepest descent doesn't point to the minimum



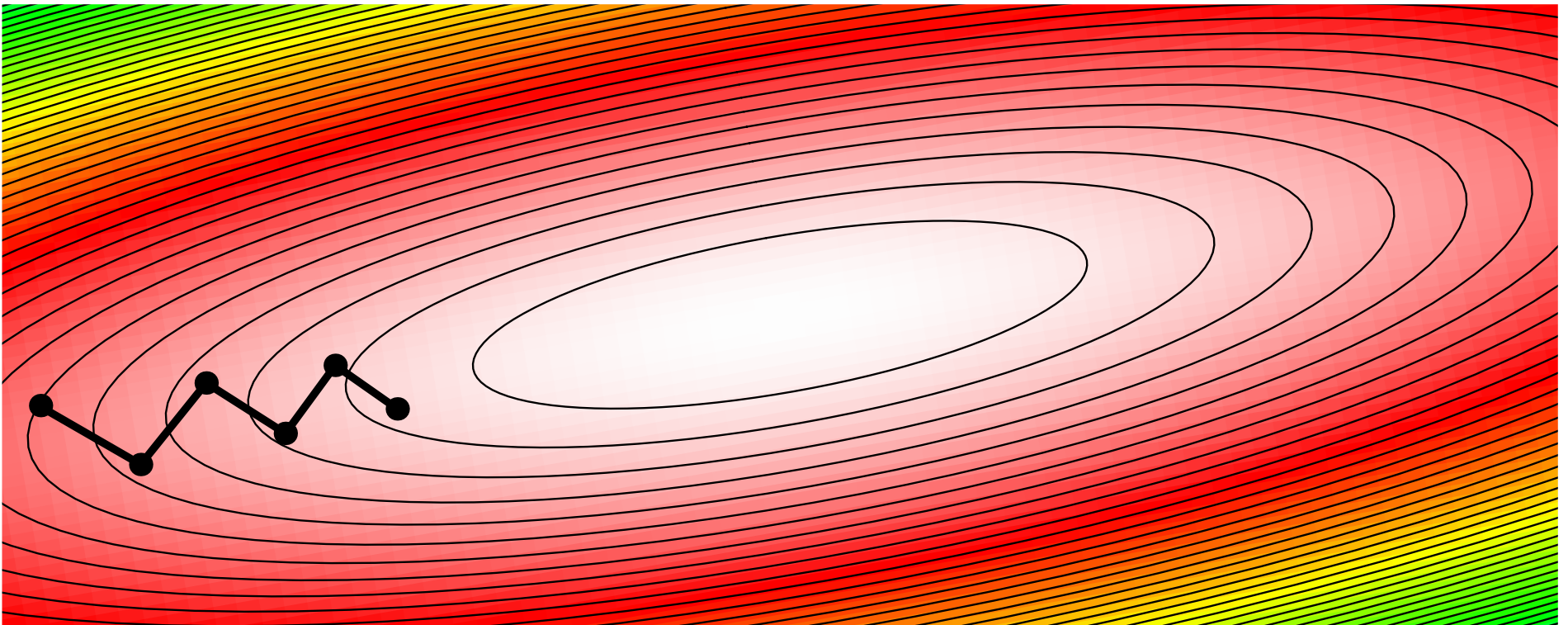
# Getting There Quicker

- Increasing the step size speeds up convergence, but the direction of steepest descent doesn't point to the minimum



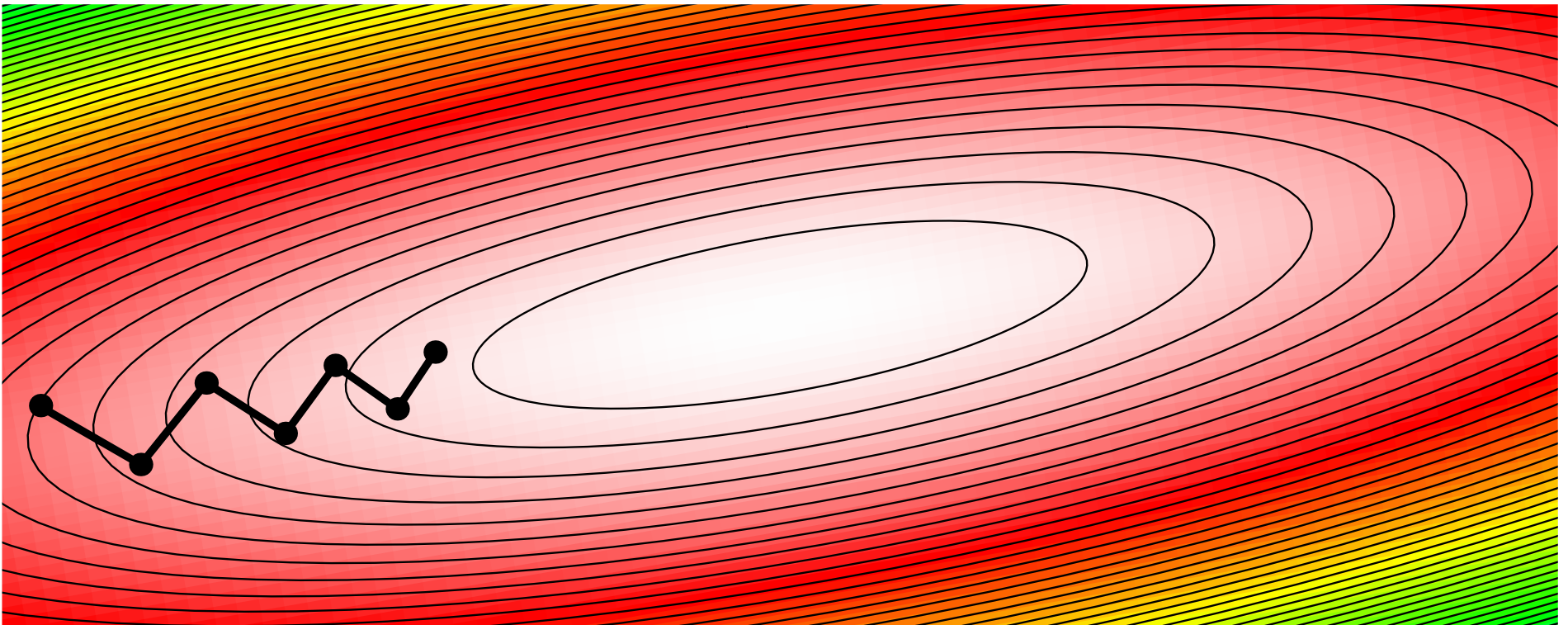
# Getting There Quicker

- Increasing the step size speeds up convergence, but the direction of steepest descent doesn't point to the minimum



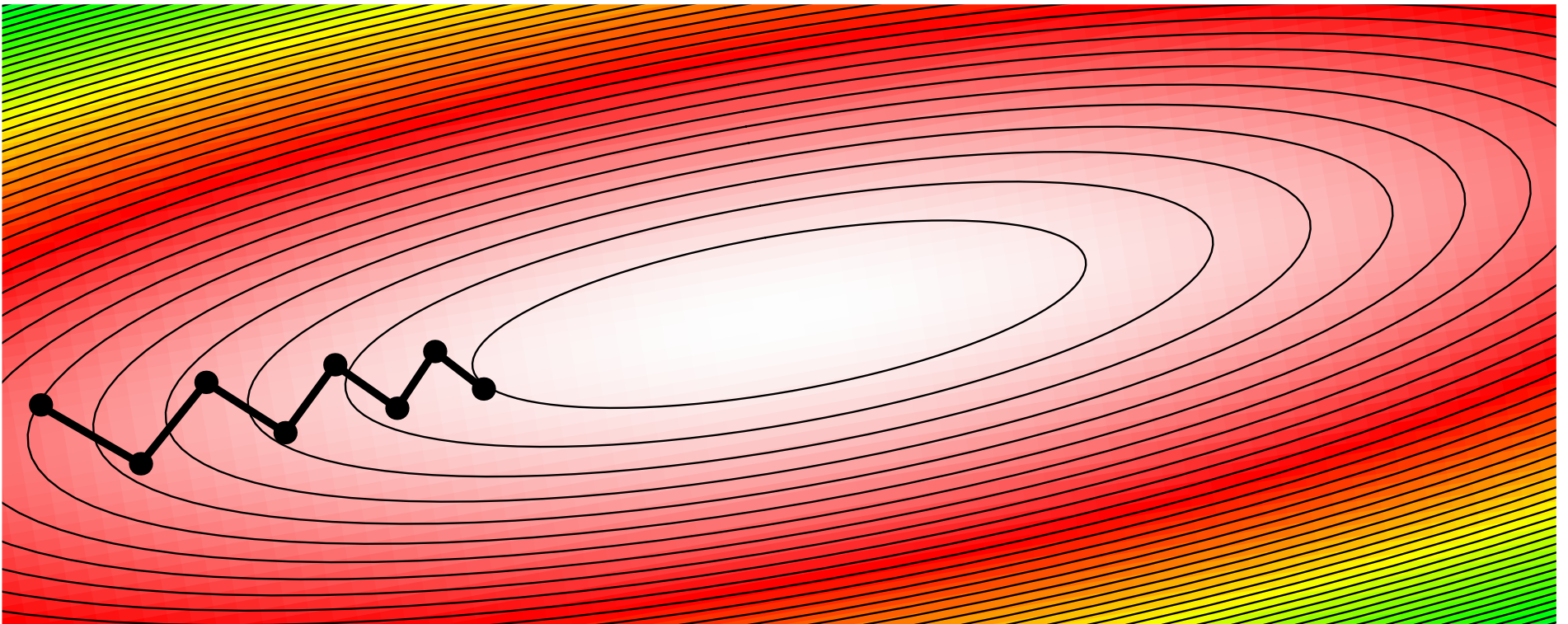
# Getting There Quicker

- Increasing the step size speeds up convergence, but the direction of steepest descent doesn't point to the minimum



# Getting There Quicker

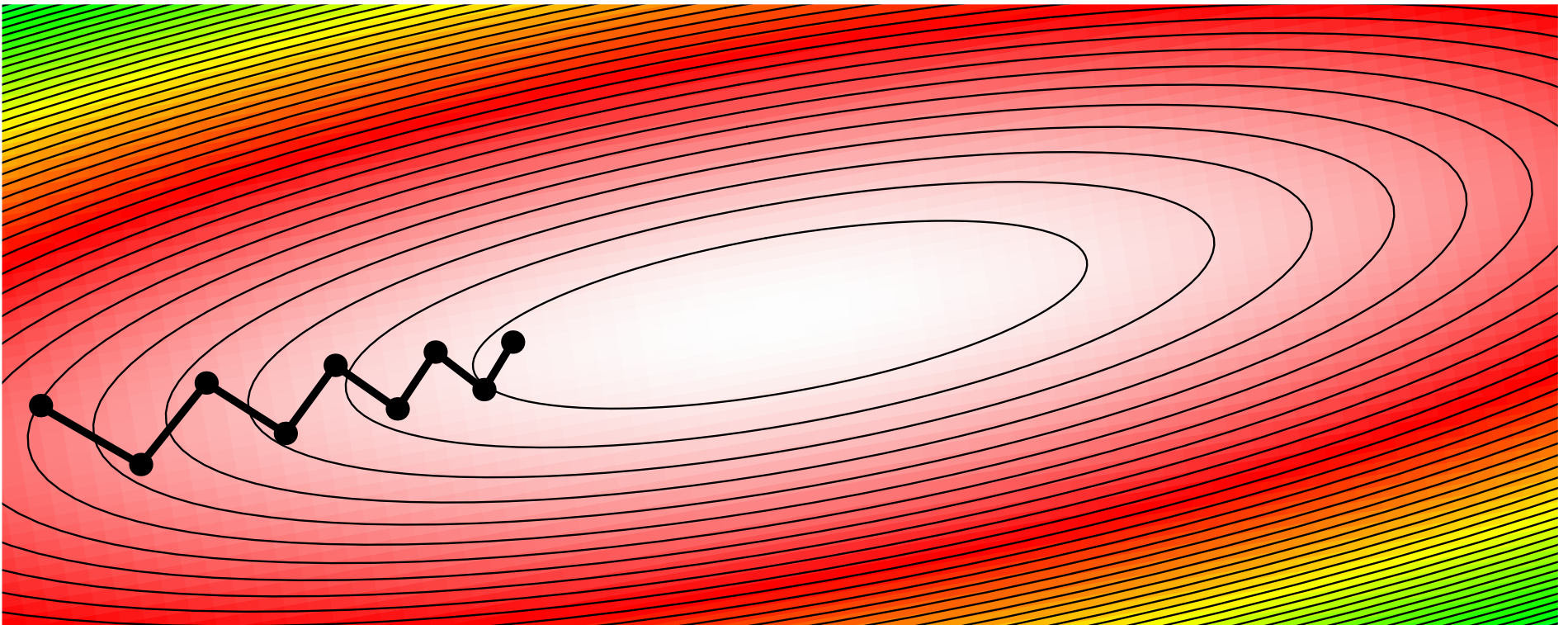
- Increasing the step size speeds up convergence, but the direction of steepest descent doesn't point to the minimum





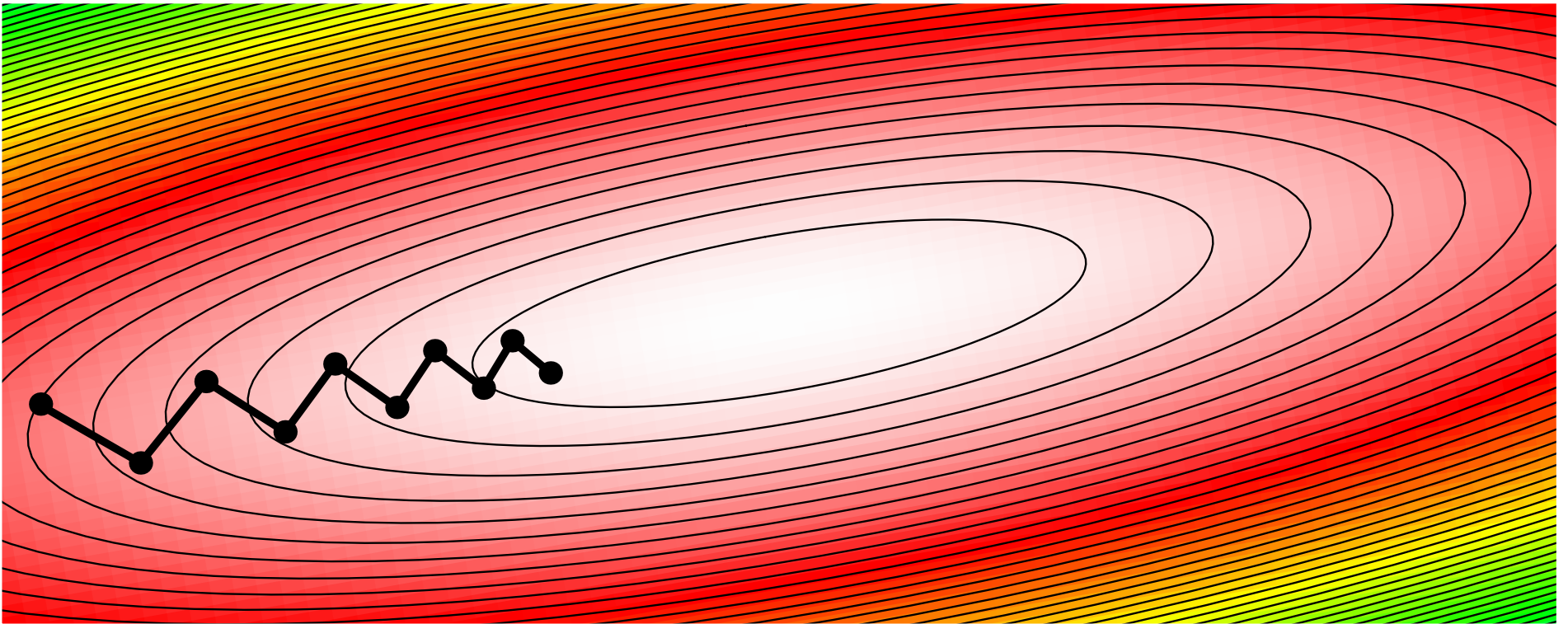
# Getting There Quicker

- Increasing the step size speeds up convergence, but the direction of steepest descent doesn't point to the minimum



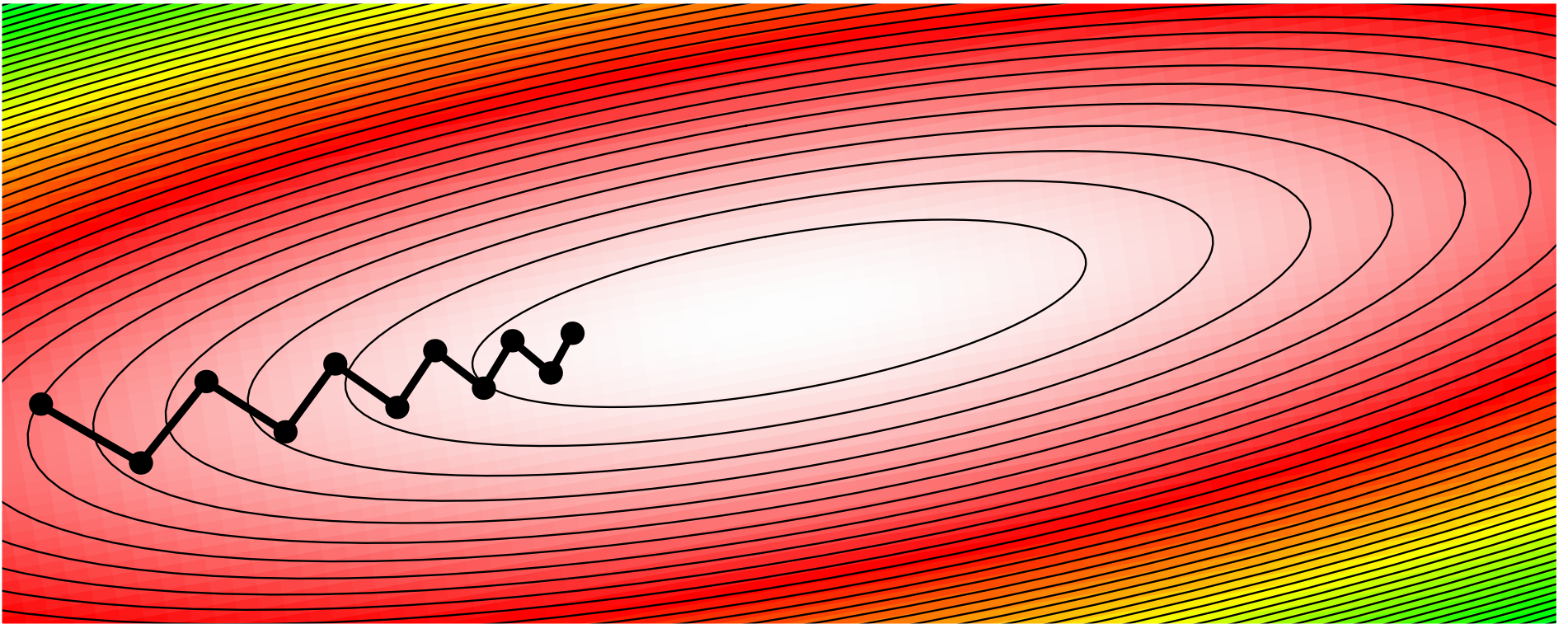
# Getting There Quicker

- Increasing the step size speeds up convergence, but the direction of steepest descent doesn't point to the minimum



# Getting There Quicker

- Increasing the step size speeds up convergence, but the direction of steepest descent doesn't point to the minimum



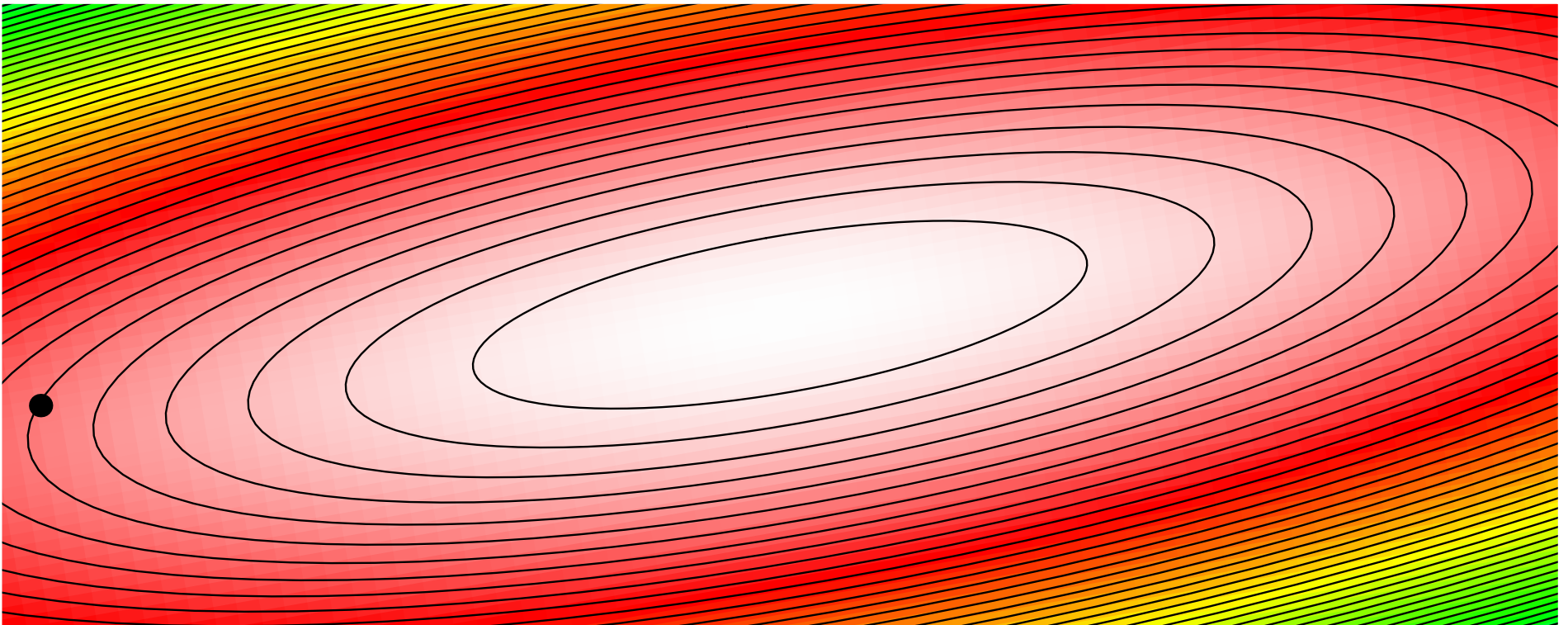


# More Haste Less Speed

- Increasing the step size, just a little further, increases the rate of converge in one direction, but . . .

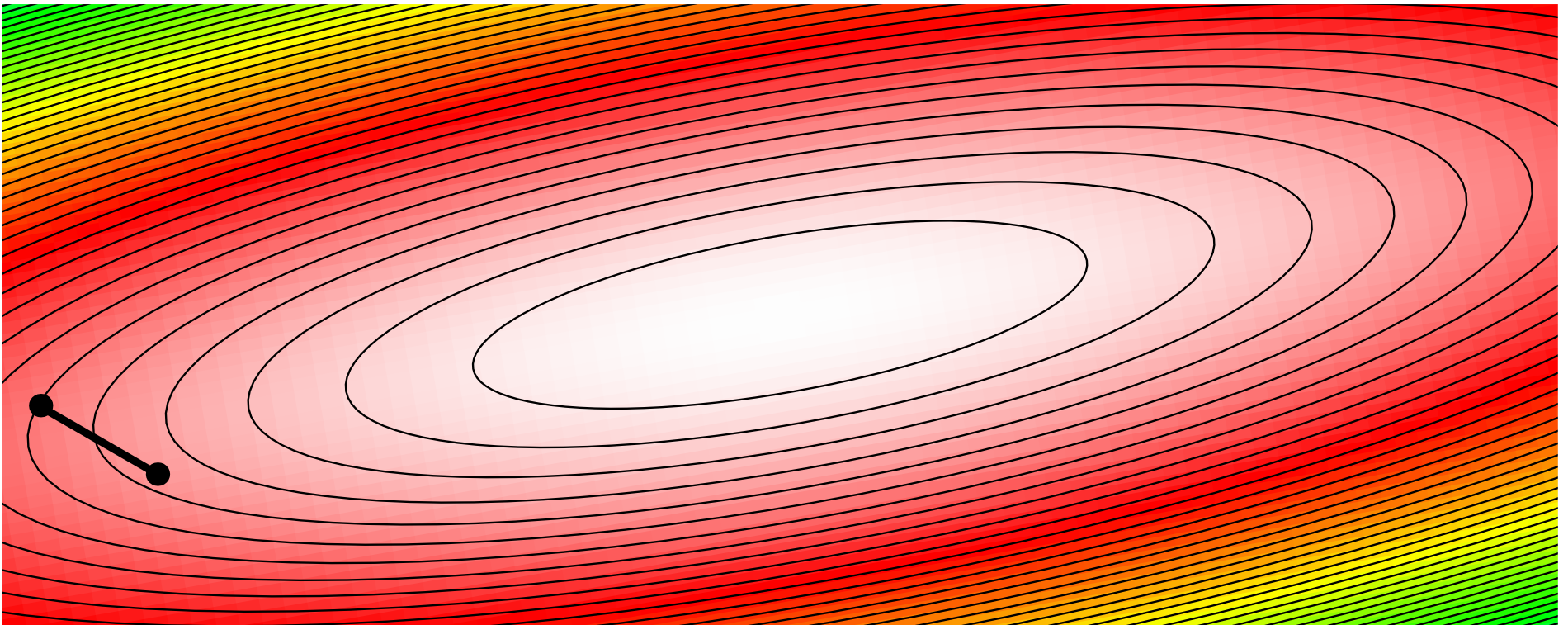
# More Haste Less Speed

- Increasing the step size, just a little further, increases the rate of converge in one direction, but . . .



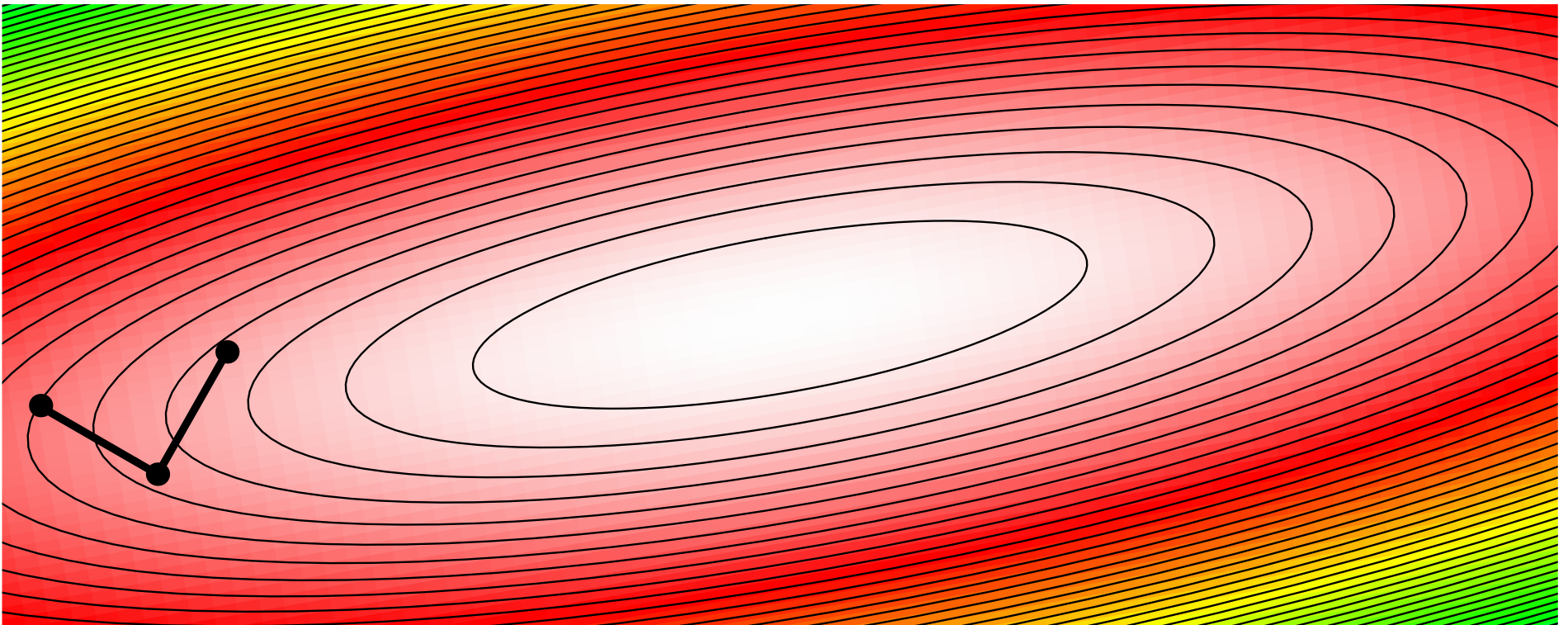
# More Haste Less Speed

- Increasing the step size, just a little further, increases the rate of converge in one direction, but . . .



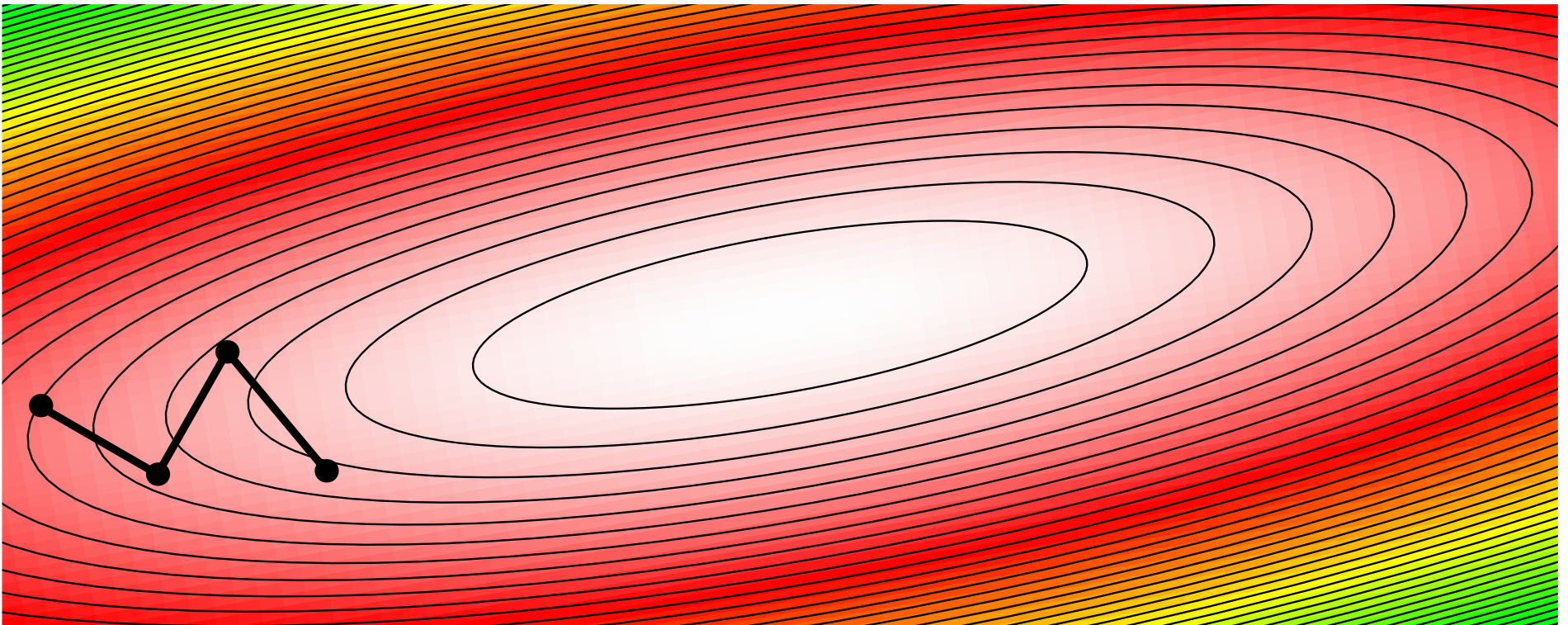
# More Haste Less Speed

- Increasing the step size, just a little further, increases the rate of converge in one direction, but . . .



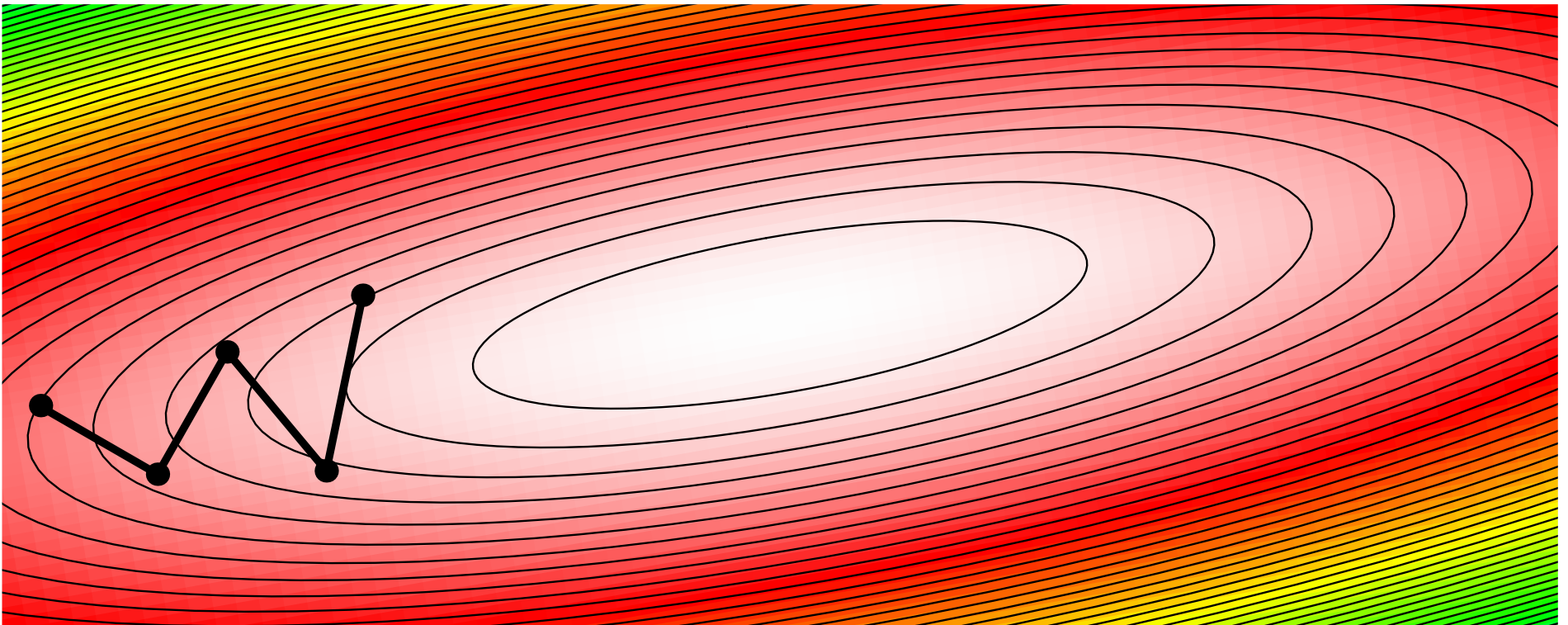
# More Haste Less Speed

- Increasing the step size, just a little further, increases the rate of converge in one direction, but . . .



# More Haste Less Speed

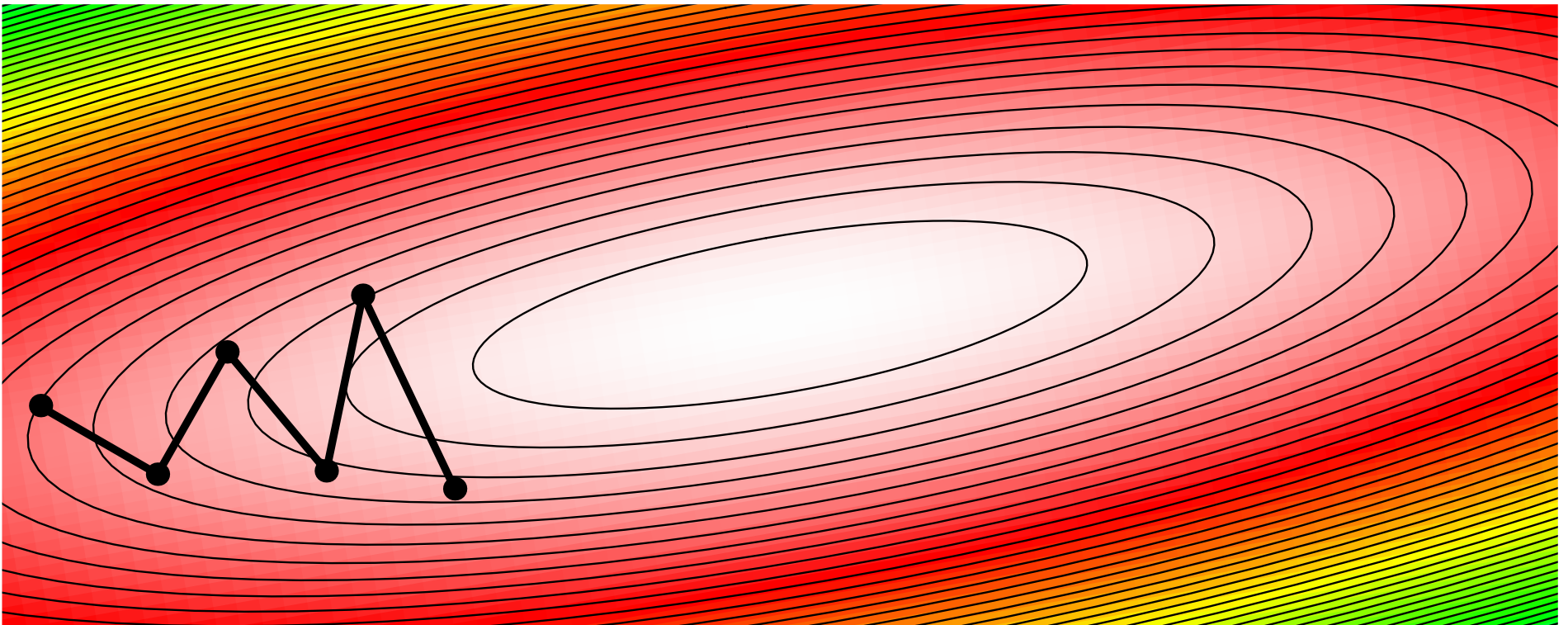
- Increasing the step size, just a little further, increases the rate of converge in one direction, but . . .





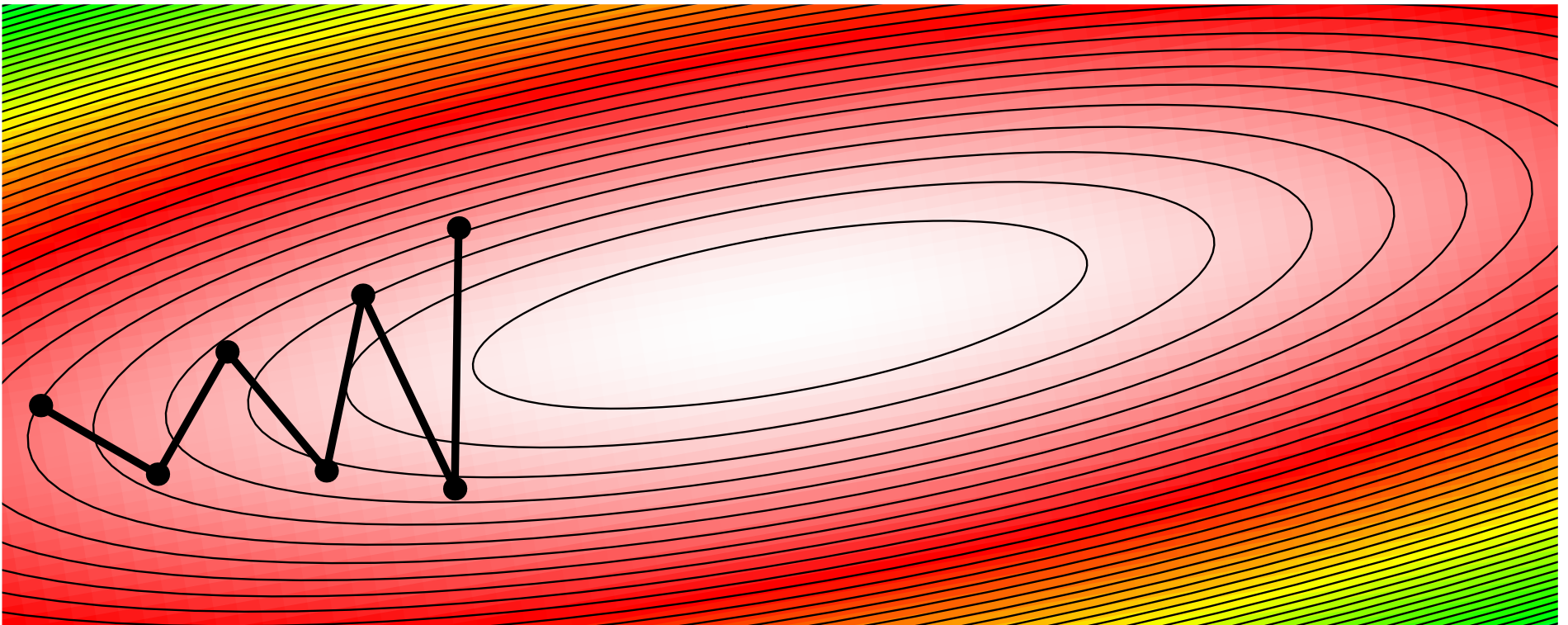
# More Haste Less Speed

- Increasing the step size, just a little further, increases the rate of converge in one direction, but . . .



# More Haste Less Speed

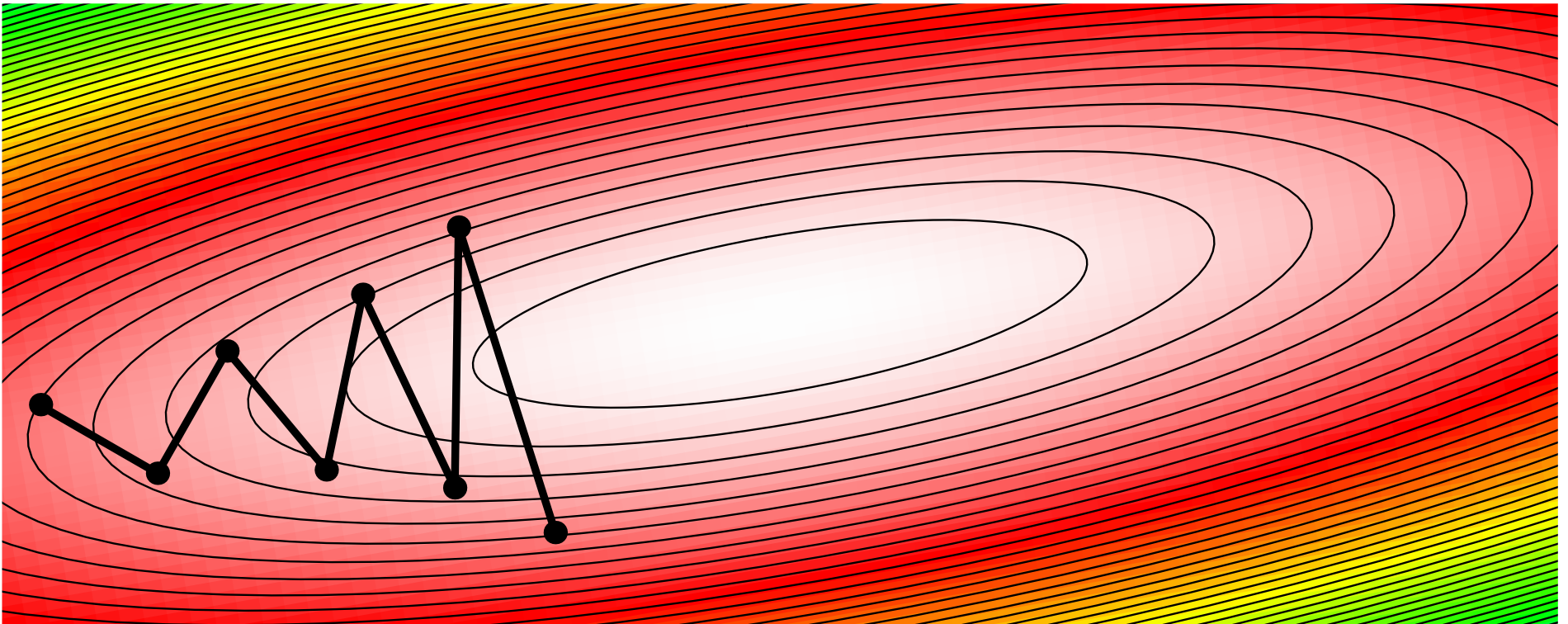
- Increasing the step size, just a little further, increases the rate of converge in one direction, but . . .





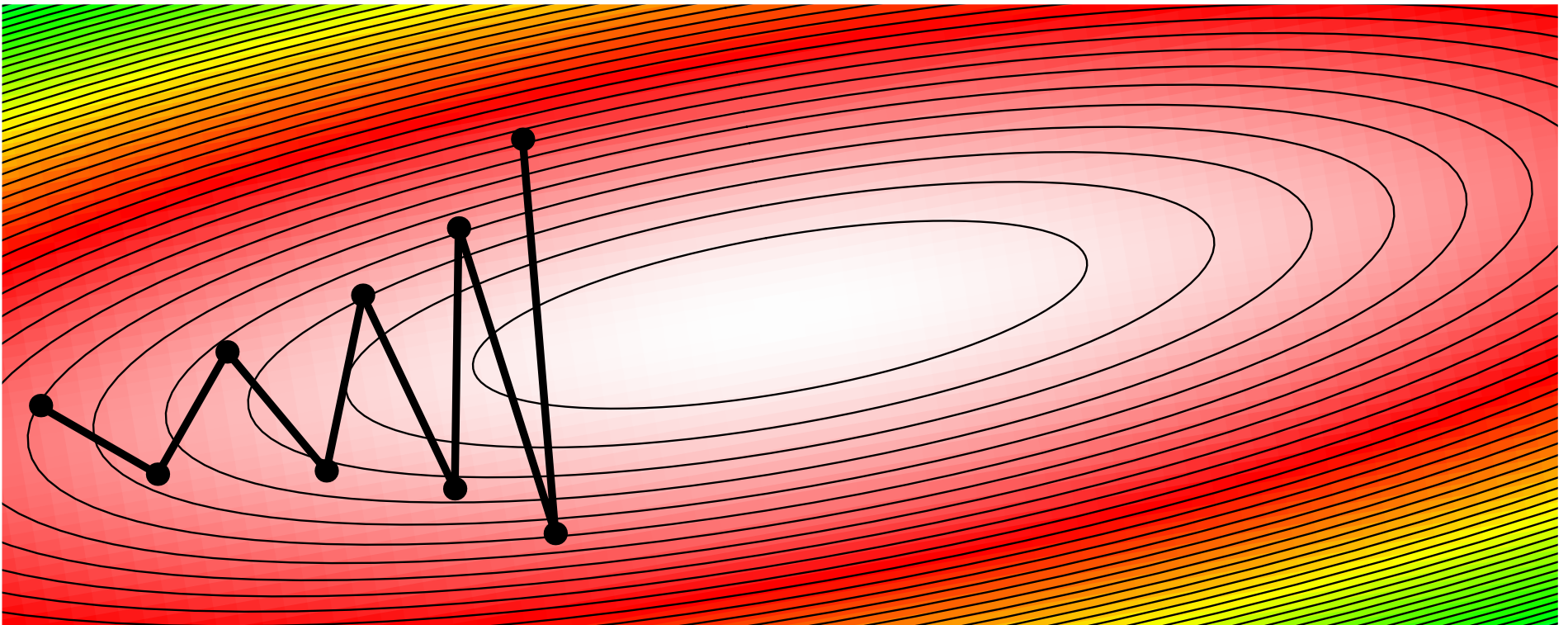
# More Haste Less Speed

- Increasing the step size, just a little further, increases the rate of converge in one direction, but . . .



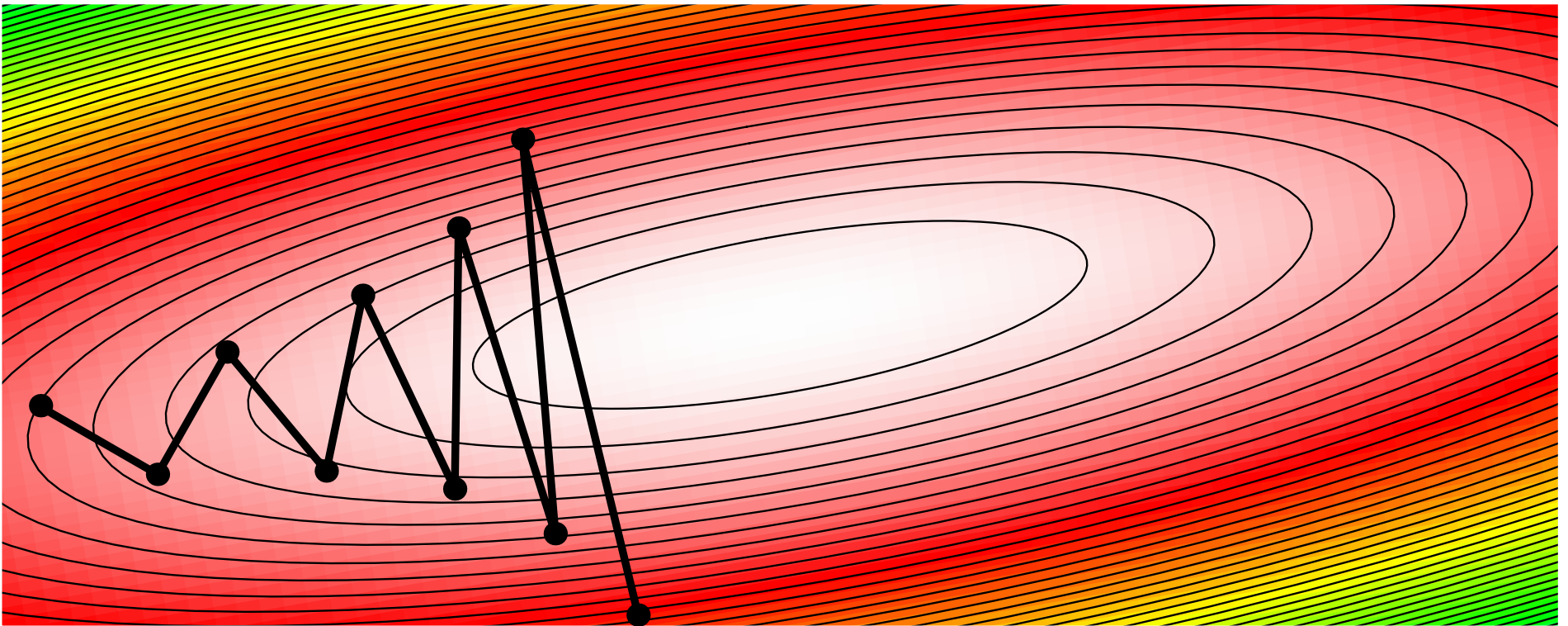
# More Haste Less Speed

- Increasing the step size, just a little further, increases the rate of converge in one direction, but . . .



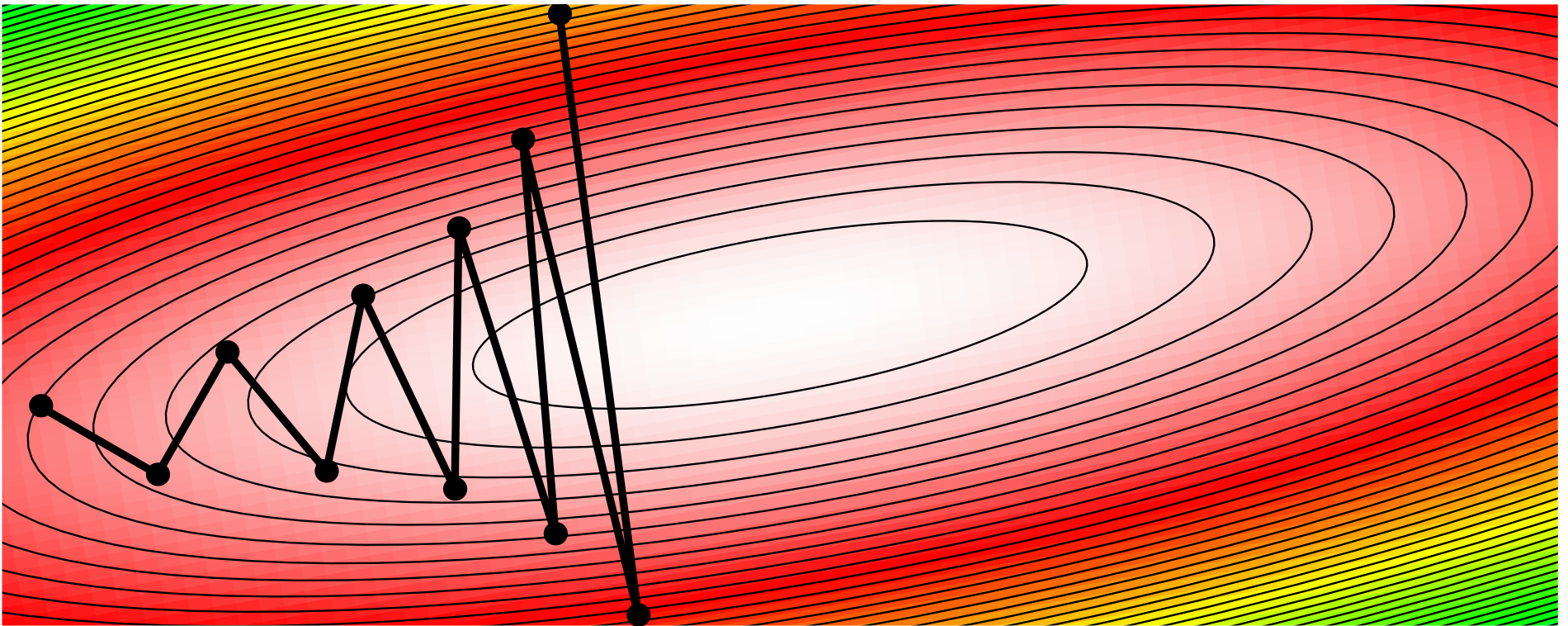
# More Haste Less Speed

- Increasing the step size, just a little further, increases the rate of converge in one direction, but . . .



# More Haste Less Speed

- Increasing the step size, just a little further, increases the rate of converge in one direction, but . . .

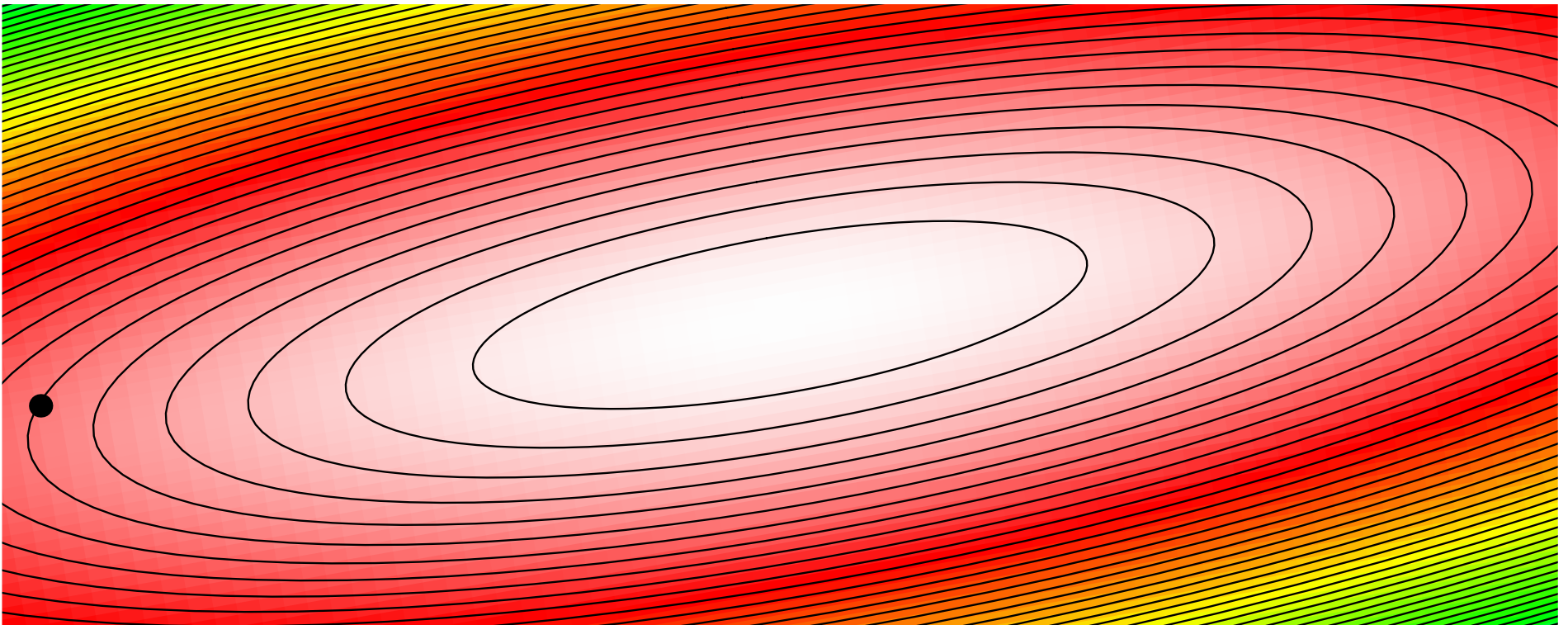


# Line Minimisation

- We can systematically seek the minimum along a line of the gradient

# Line Minimisation

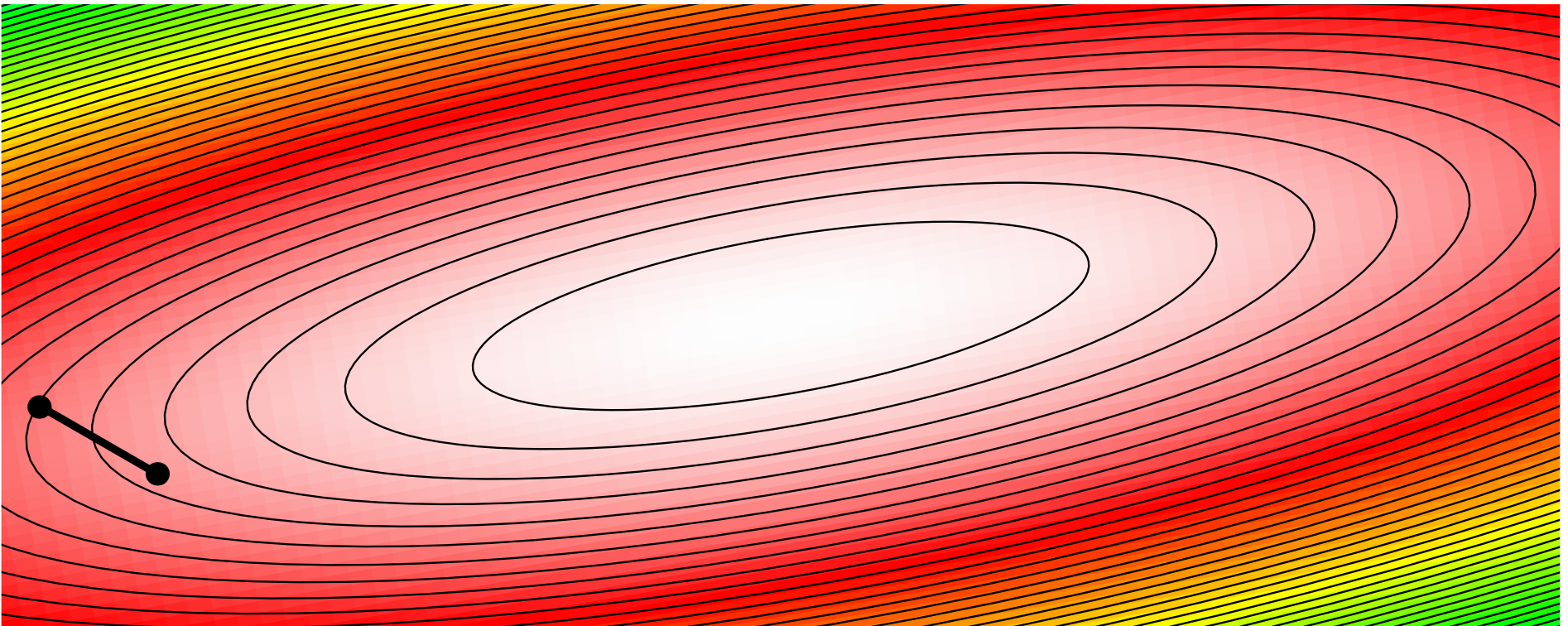
- We can systematically seek the minimum along a line of the gradient





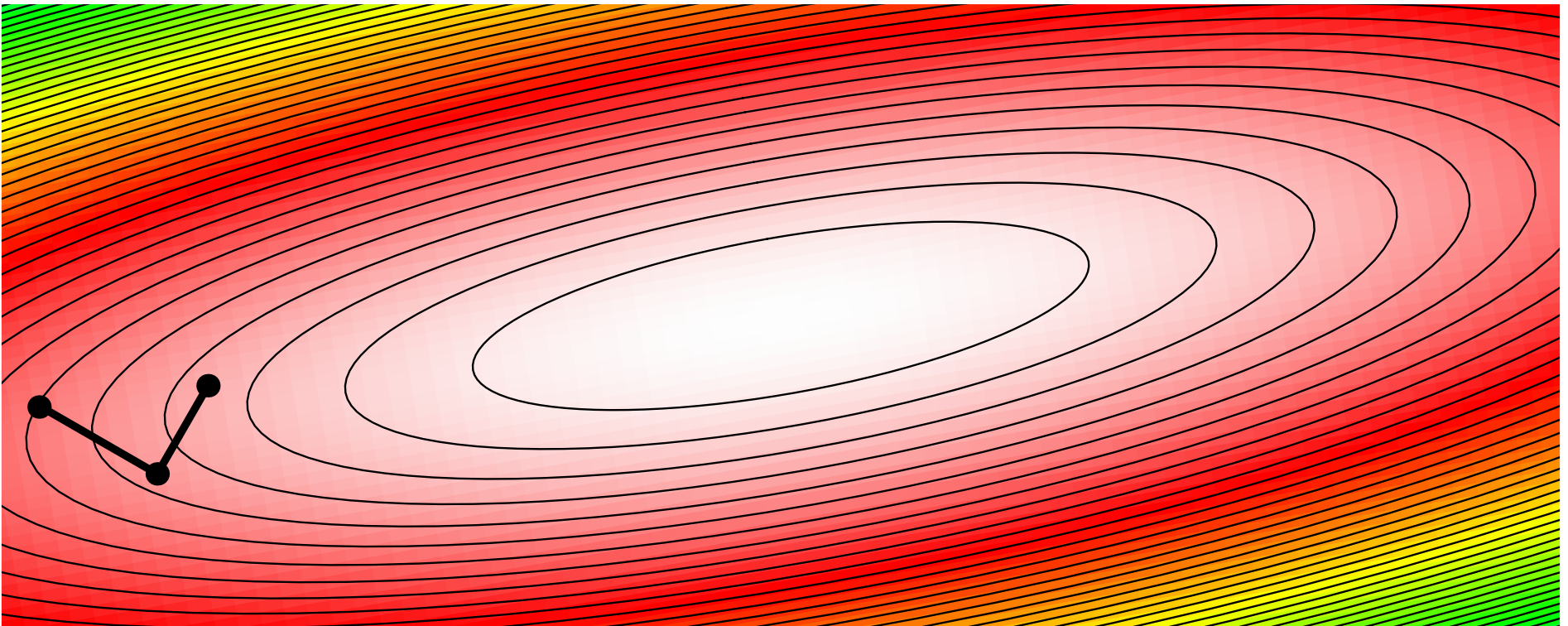
# Line Minimisation

- We can systematically seek the minimum along a line of the gradient



# Line Minimisation

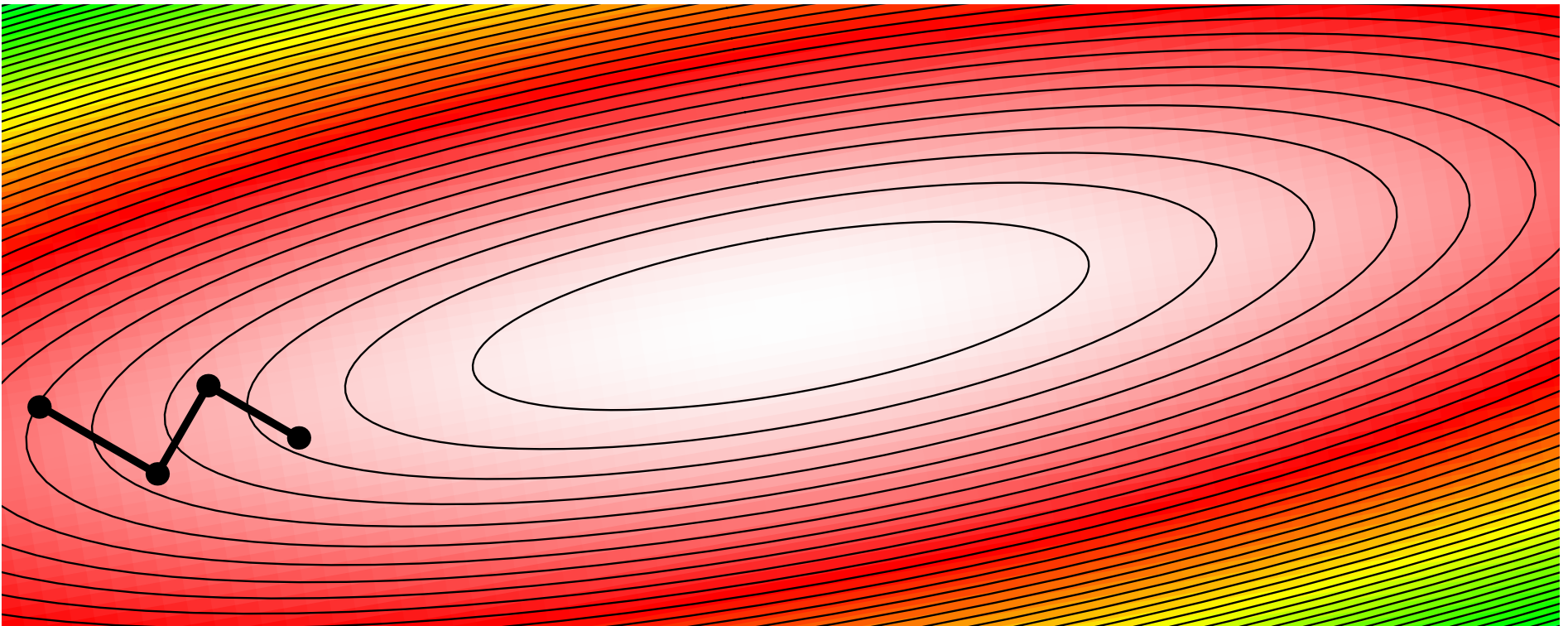
- We can systematically seek the minimum along a line of the gradient





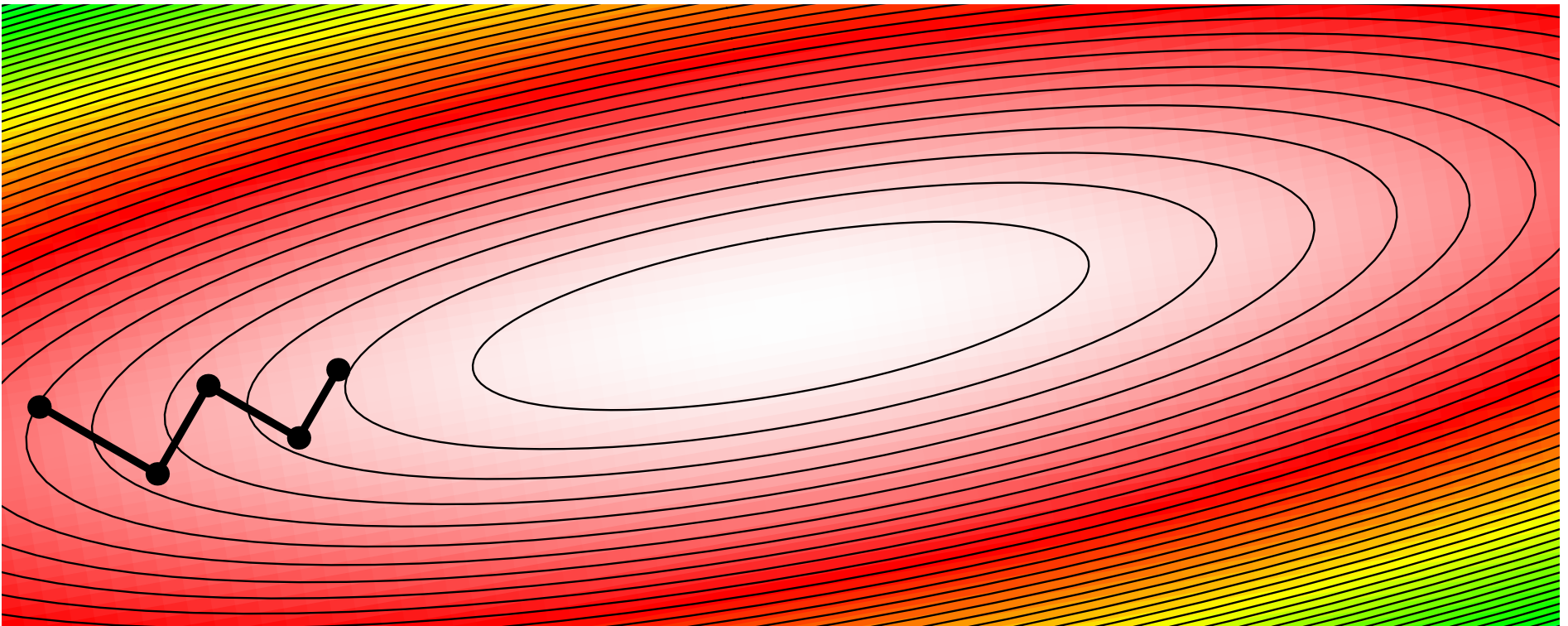
# Line Minimisation

- We can systematically seek the minimum along a line of the gradient



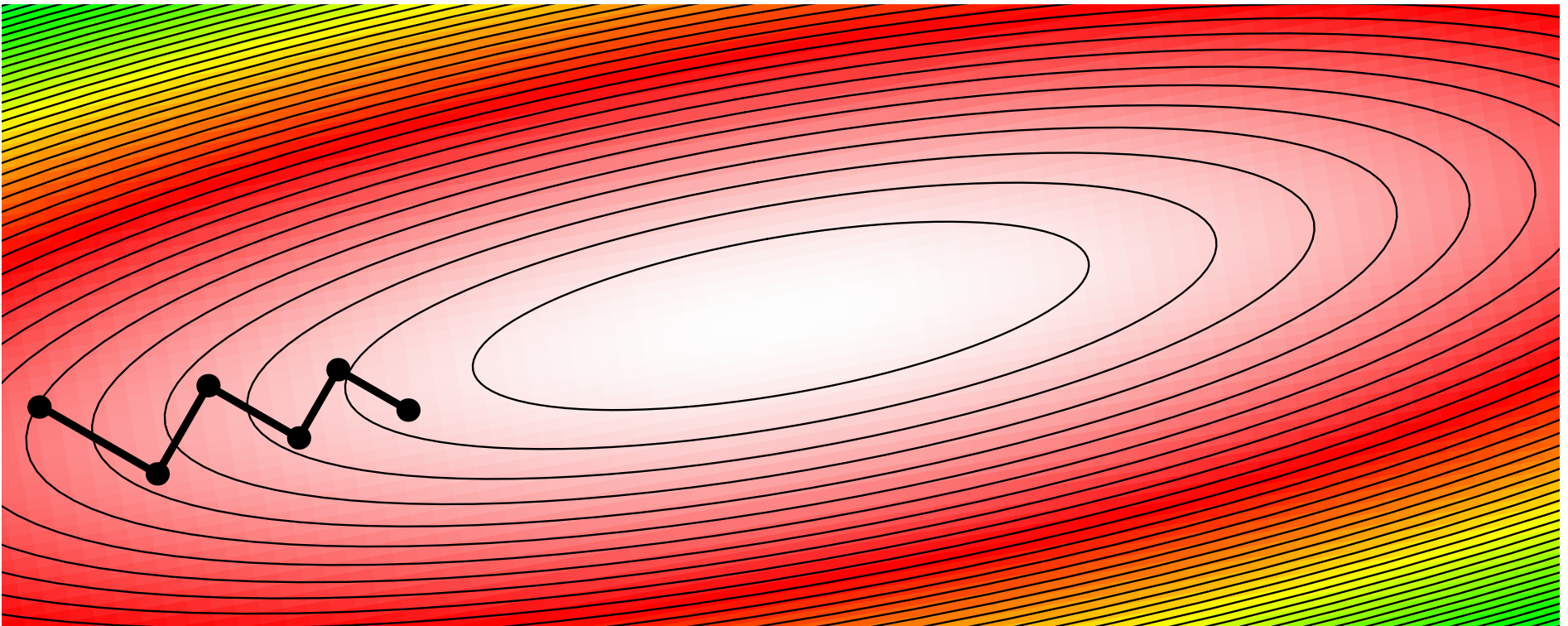
# Line Minimisation

- We can systematically seek the minimum along a line of the gradient



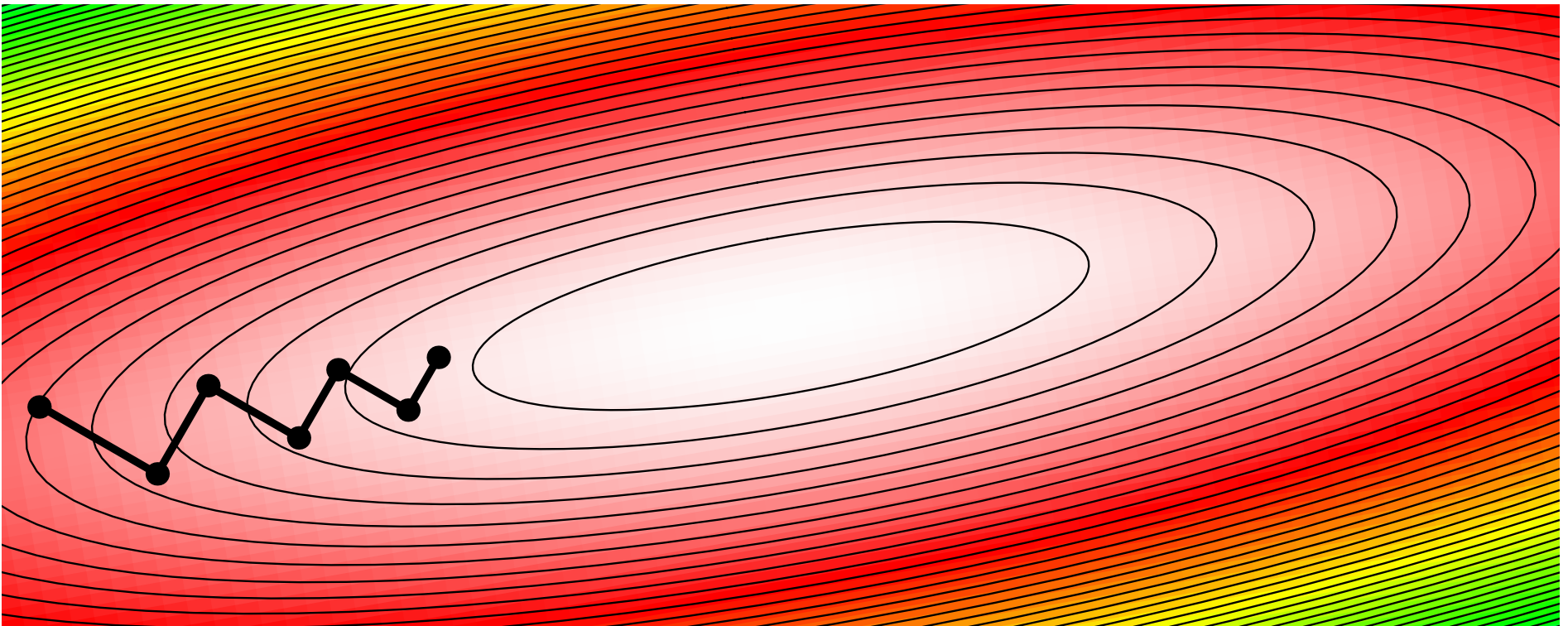
# Line Minimisation

- We can systematically seek the minimum along a line of the gradient



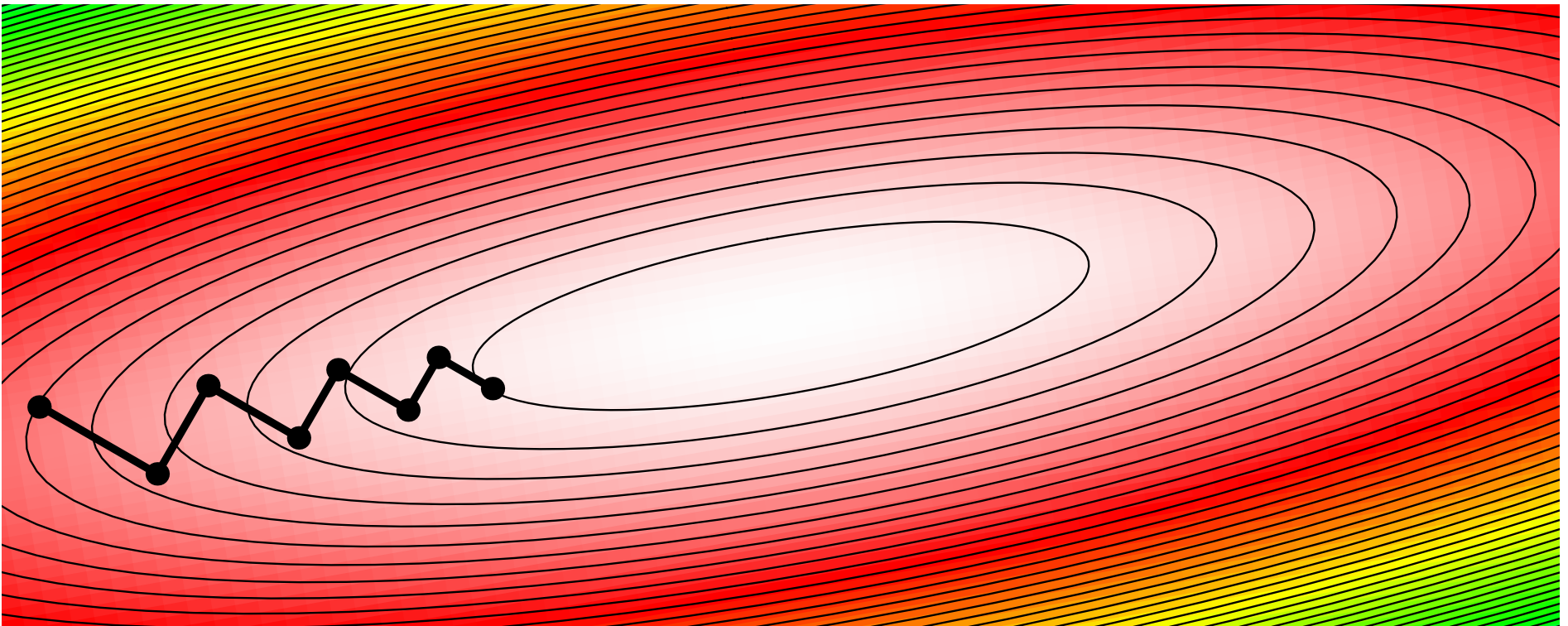
# Line Minimisation

- We can systematically seek the minimum along a line of the gradient



# Line Minimisation

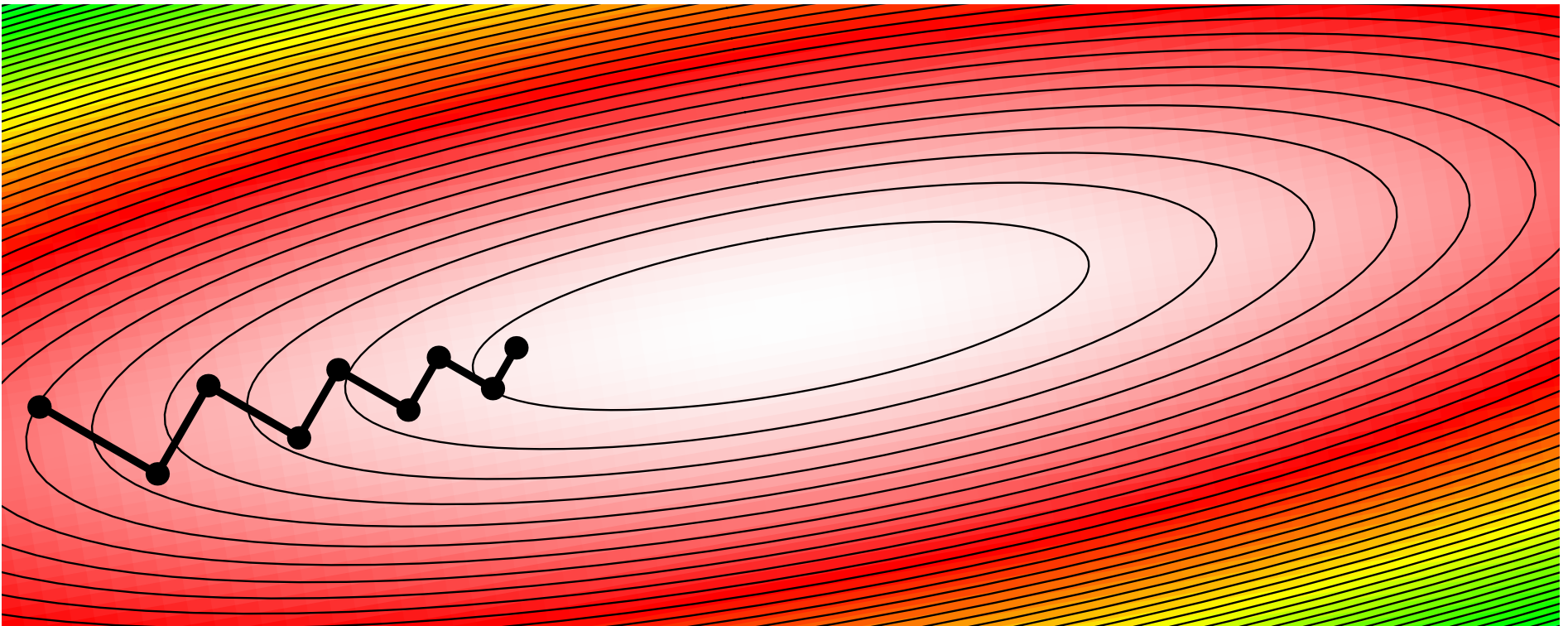
- We can systematically seek the minimum along a line of the gradient





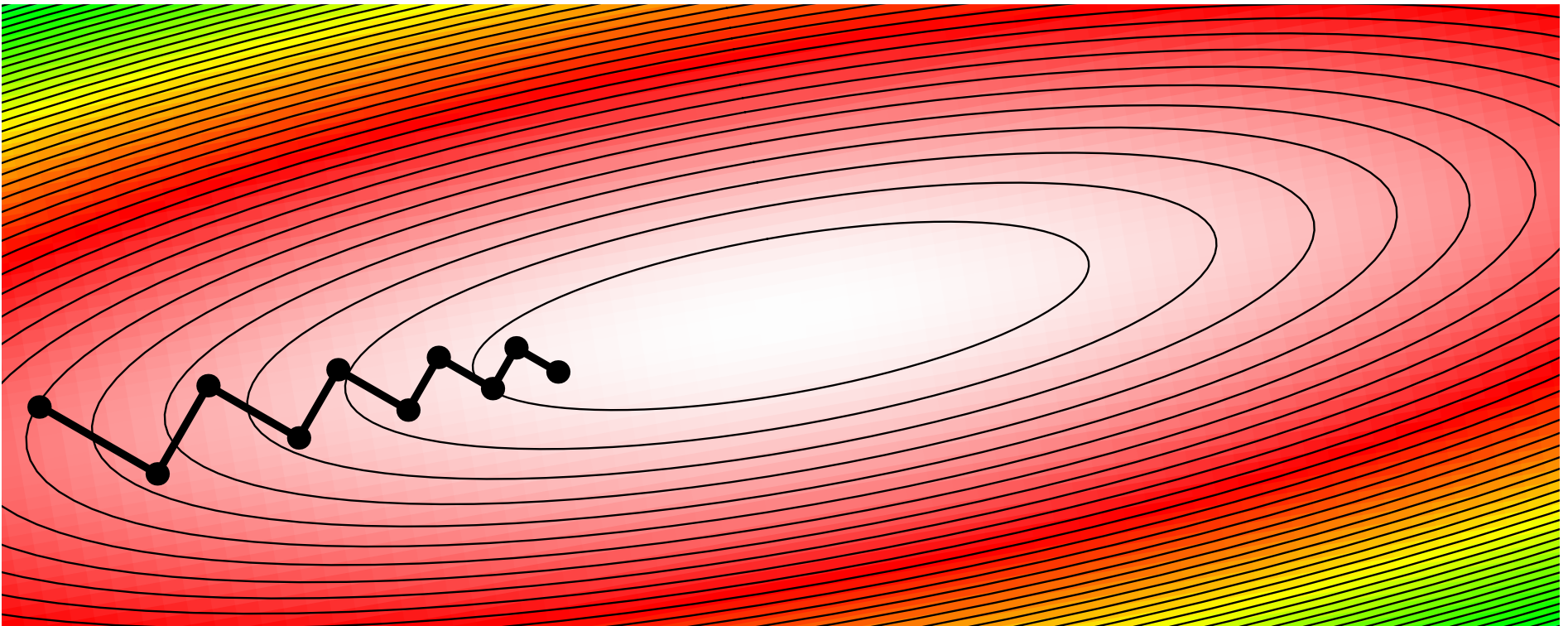
# Line Minimisation

- We can systematically seek the minimum along a line of the gradient



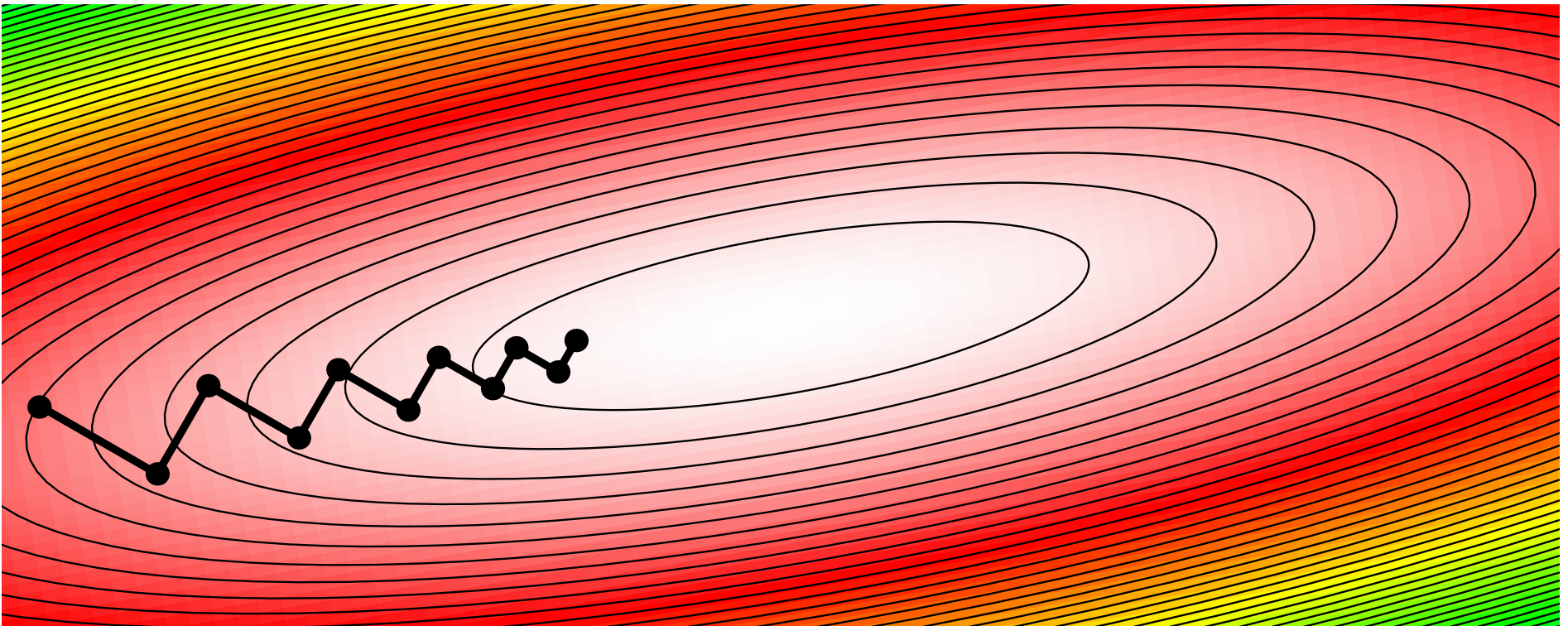
# Line Minimisation

- We can systematically seek the minimum along a line of the gradient



# Line Minimisation

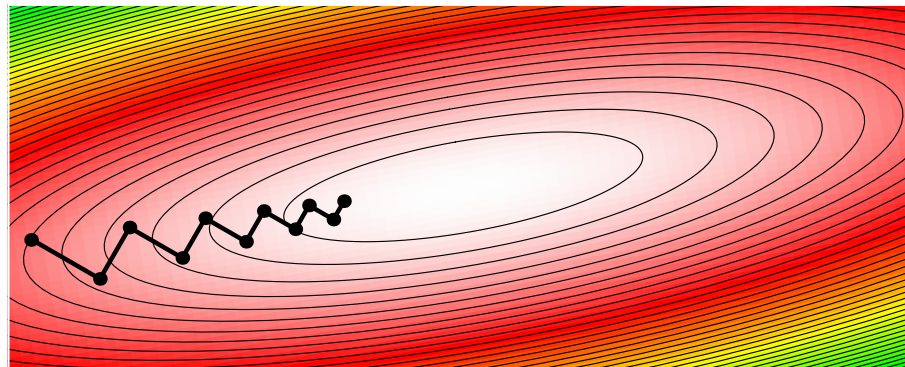
- We can systematically seek the minimum along a line of the gradient





# Zig-Zag

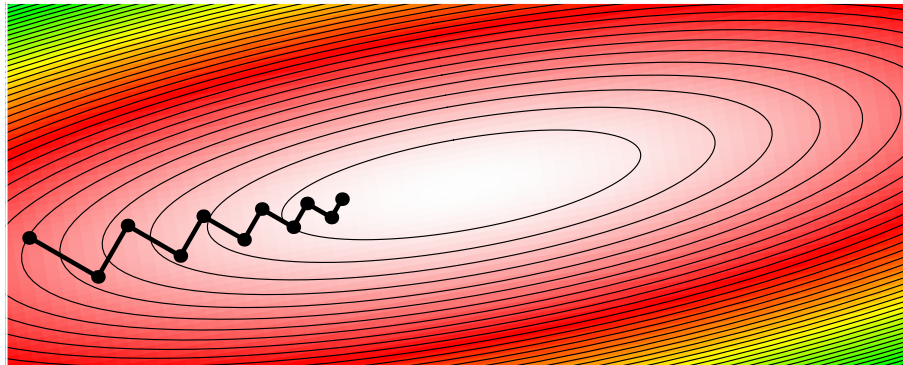
- Note that in high dimensions gradient descent tends to zigzag



- If we computed the Hessian and used Newton's method we would jump straight to the minimum if we were in a quadratic potential
- However computing the Hessian is time consuming and misleading if we are not in a quadratic potential (i.e. far from the optimum)

# Zig-Zag

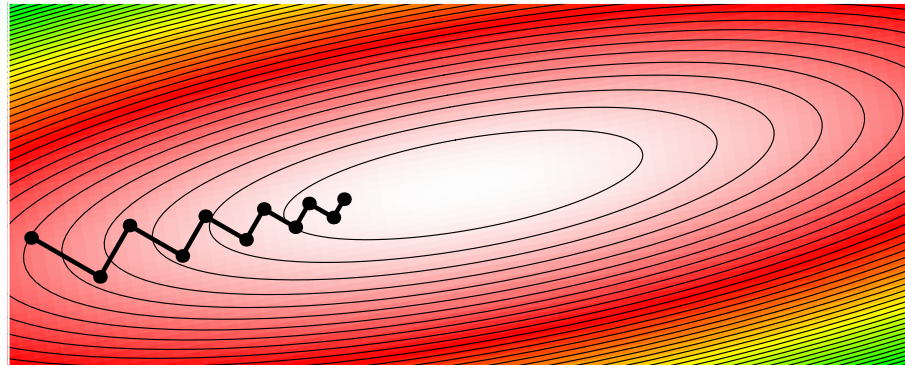
- Note that in high dimensions gradient descent tends to zigzag



- If we computed the Hessian and used Newton's method we would jump straight to the minimum if we were in a quadratic potential
- However computing the Hessian is time consuming and misleading if we are not in a quadratic potential (i.e. far from the optimum)

# Zig-Zag

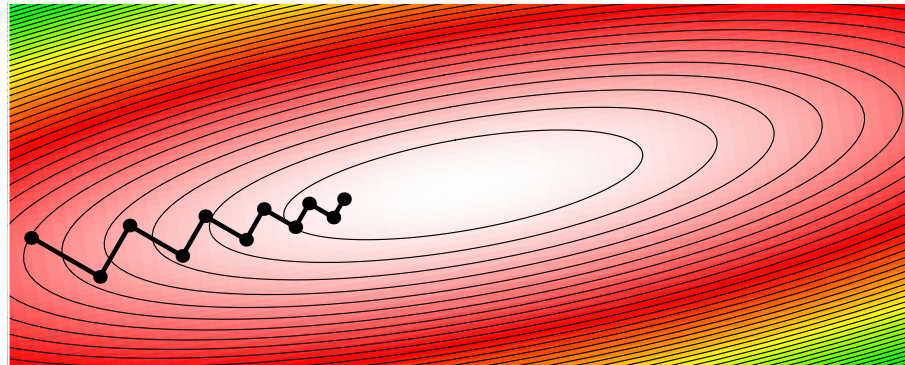
- Note that in high dimensions gradient descent tends to zigzag



- If we computed the Hessian and used Newton's method we would jump straight to the minimum if we were in a quadratic potential
- However computing the Hessian is time consuming and misleading if we are not in a quadratic potential (i.e. far from the optimum)

# Zig-Zag

- Note that in high dimensions gradient descent tends to zigzag



- If we computed the Hessian and used Newton's method we would jump straight to the minimum if we were in a quadratic potential
- However computing the Hessian is time consuming and misleading if we are not in a quadratic potential (i.e. far from the optimum)

# Best Optimisation Algorithms

- The best optimisation algorithms compute an approximation of the Hessian
- E.g. Conjugate gradient
  - ★ Performs Line Minimisation
  - ★ Uses gradient, but does not go along it
  - ★ Reaches quadratic minimum in  $N$  steps
- E.g. Levenberg-Marquardt
  - ★ Used on least squares problem only
  - ★ Uses linear approximation of function to approximate Hessian
  - ★ Adapts from hill-climbing to Newton method
  - ★ Avoids line-minimisation

# Best Optimisation Algorithms

- The best optimisation algorithms compute an approximation of the Hessian
- E.g. Conjugate gradient
  - ★ Performs Line Minimisation
  - ★ Uses gradient, but does not go along it
  - ★ Reaches quadratic minimum in  $N$  steps
- E.g. Levenberg-Marquardt
  - ★ Used on least squares problem only
  - ★ Uses linear approximation of function to approximate Hessian
  - ★ Adapts from hill-climbing to Newton method
  - ★ Avoids line-minimisation

# Best Optimisation Algorithms

- The best optimisation algorithms compute an approximation of the Hessian
- E.g. Conjugate gradient
  - ★ Performs Line Minimisation
  - ★ Uses gradient, but does not go along it
  - ★ Reaches quadratic minimum in  $N$  steps
- E.g. Levenberg-Marquardt
  - ★ Used on least squares problem only
  - ★ Uses linear approximation of function to approximate Hessian
  - ★ Adapts from hill-climbing to Newton method
  - ★ Avoids line-minimisation

# Best Optimisation Algorithms

- The best optimisation algorithms compute an approximation of the Hessian
- E.g. Conjugate gradient
  - ★ Performs Line Minimisation
  - ★ Uses gradient, but does not go along it
  - ★ Reaches quadratic minimum in  $N$  steps
- E.g. Levenberg-Marquardt
  - ★ Used on least squares problem only
  - ★ Uses linear approximation of function to approximate Hessian
  - ★ Adapts from hill-climbing to Newton method
  - ★ Avoids line-minimisation



# Best Optimisation Algorithms

- The best optimisation algorithms compute an approximation of the Hessian
- E.g. Conjugate gradient
  - ★ Performs Line Minimisation
  - ★ Uses gradient, but does not go along it
  - ★ Reaches quadratic minimum in  $N$  steps
- E.g. Levenberg-Marquardt
  - ★ Used on least squares problem only
  - ★ Uses linear approximation of function to approximate Hessian
  - ★ Adapts from hill-climbing to Newton method
  - ★ Avoids line-minimisation

# Best Optimisation Algorithms

- The best optimisation algorithms compute an approximation of the Hessian
- E.g. Conjugate gradient
  - ★ Performs Line Minimisation
  - ★ Uses gradient, but does not go along it
  - ★ Reaches quadratic minimum in  $N$  steps
- E.g. Levenberg-Marquardt
  - ★ Used on least squares problem only
  - ★ Uses linear approximation of function to approximate Hessian
  - ★ Adapts from hill-climbing to Newton method
  - ★ Avoids line-minimisation

# Best Optimisation Algorithms

- The best optimisation algorithms compute an approximation of the Hessian
- E.g. Conjugate gradient
  - ★ Performs Line Minimisation
  - ★ Uses gradient, but does not go along it
  - ★ Reaches quadratic minimum in  $N$  steps
- E.g. Levenberg-Marquardt
  - ★ Used on least squares problem only
  - ★ Uses linear approximation of function to approximate Hessian
  - ★ Adapts from hill-climbing to Newton method
  - ★ Avoids line-minimisation

# Best Optimisation Algorithms

- The best optimisation algorithms compute an approximation of the Hessian
- E.g. Conjugate gradient
  - ★ Performs Line Minimisation
  - ★ Uses gradient, but does not go along it
  - ★ Reaches quadratic minimum in  $N$  steps
- E.g. Levenberg-Marquardt
  - ★ Used on least squares problem only
  - ★ Uses linear approximation of function to approximate Hessian
  - ★ Adapts from hill-climbing to Newton method
  - ★ Avoids line-minimisation

# Best Optimisation Algorithms

- The best optimisation algorithms compute an approximation of the Hessian
- E.g. Conjugate gradient
  - ★ Performs Line Minimisation
  - ★ Uses gradient, but does not go along it
  - ★ Reaches quadratic minimum in  $N$  steps
- E.g. Levenberg-Marquardt
  - ★ Used on least squares problem only
  - ★ Uses linear approximation of function to approximate Hessian
  - ★ Adapts from hill-climbing to Newton method
  - ★ Avoids line-minimisation

# Best Optimisation Algorithms

- The best optimisation algorithms compute an approximation of the Hessian
- E.g. Conjugate gradient
  - ★ Performs Line Minimisation
  - ★ Uses gradient, but does not go along it
  - ★ Reaches quadratic minimum in  $N$  steps
- E.g. Levenberg-Marquardt
  - ★ Used on least squares problem only
  - ★ Uses linear approximation of function to approximate Hessian
  - ★ Adapts from hill-climbing to Newton method
  - ★ **Avoids line-minimisation**

# Levenberg-Marquardt

- Want to minimise  $\|\epsilon(\mathbf{w})\|^2$  where  $\epsilon_i(\mathbf{w}) = f(\mathbf{x}_i|\mathbf{w}) - y_i$
- Use linear approximation

$$\epsilon_i(\mathbf{w}) \approx \epsilon_i(\mathbf{w}^{(k)}) + (\mathbf{w} - \mathbf{w}^{(k)})^\top \nabla \epsilon_i(\mathbf{w}^{(k)})$$

with  $\nabla \epsilon_i(\mathbf{w}^{(k)}) = \nabla f(\mathbf{x}_i|\mathbf{w}^{(k)})$

- Solve quadratic minimisation of approximate error  
 $\operatorname{argmin}_{\mathbf{w}} L_{\text{approx}}(\mathbf{w})$  with  $\mathbf{J} = \nabla \epsilon(\mathbf{w}^{(k)})$

$$\begin{aligned} L_{\text{approx}}(\mathbf{w}) &= \|\epsilon(\mathbf{w}^{(k)}) + \mathbf{J}(\mathbf{w} - \mathbf{w}^{(k)})\|^2 \\ &= \epsilon(\mathbf{w}^{(k)})^\top \epsilon(\mathbf{w}^{(k)}) + 2(\mathbf{w} - \mathbf{w}^{(k)})^\top \mathbf{J}^\top \epsilon(\mathbf{w}^{(k)}) \\ &\quad + (\mathbf{w} - \mathbf{w}^{(k)})^\top \mathbf{J}^\top \mathbf{J} (\mathbf{w} - \mathbf{w}^{(k)}) \end{aligned}$$

# Levenberg-Marquardt

- Want to minimise  $\|\epsilon(\mathbf{w})\|^2$  where  $\epsilon_i(\mathbf{w}) = f(\mathbf{x}_i|\mathbf{w}) - y_i$
- Use linear approximation

$$\epsilon_i(\mathbf{w}) \approx \epsilon_i(\mathbf{w}^{(k)}) + (\mathbf{w} - \mathbf{w}^{(k)})^\top \nabla \epsilon_i(\mathbf{w}^{(k)})$$

$$\text{with } \nabla \epsilon_i(\mathbf{w}^{(k)}) = \nabla f(\mathbf{x}_i|\mathbf{w}^{(k)})$$

- Solve quadratic minimisation of approximate error  
 $\operatorname{argmin}_{\mathbf{w}} L_{\text{approx}}(\mathbf{w})$  with  $\mathbf{J} = \nabla \epsilon(\mathbf{w}^{(k)})$

$$\begin{aligned} L_{\text{approx}}(\mathbf{w}) &= \|\epsilon(\mathbf{w}^{(k)}) + \mathbf{J}(\mathbf{w} - \mathbf{w}^{(k)})\|^2 \\ &= \epsilon(\mathbf{w}^{(k)})^\top \epsilon(\mathbf{w}^{(k)}) + 2(\mathbf{w} - \mathbf{w}^{(k)})^\top \mathbf{J}^\top \epsilon(\mathbf{w}^{(k)}) \\ &\quad + (\mathbf{w} - \mathbf{w}^{(k)})^\top \mathbf{J}^\top \mathbf{J} (\mathbf{w} - \mathbf{w}^{(k)}) \end{aligned}$$



# Levenberg-Marquardt

- Want to minimise  $\|\epsilon(\mathbf{w})\|^2$  where  $\epsilon_i(\mathbf{w}) = f(\mathbf{x}_i|\mathbf{w}) - y_i$
- Use linear approximation

$$\epsilon_i(\mathbf{w}) \approx \epsilon_i(\mathbf{w}^{(k)}) + (\mathbf{w} - \mathbf{w}^{(k)})^\top \nabla \epsilon_i(\mathbf{w}^{(k)})$$

with  $\nabla \epsilon_i(\mathbf{w}^{(k)}) = \nabla f(\mathbf{x}_i|\mathbf{w}^{(k)})$

- Solve quadratic minimisation of approximate error  
 $\operatorname{argmin}_{\mathbf{w}} L_{\text{approx}}(\mathbf{w})$  with  $\mathbf{J} = \nabla \epsilon(\mathbf{w}^{(k)})$

$$\begin{aligned} L_{\text{approx}}(\mathbf{w}) &= \|\epsilon(\mathbf{w}^{(k)}) + \mathbf{J}(\mathbf{w} - \mathbf{w}^{(k)})\|^2 \\ &= \epsilon(\mathbf{w}^{(k)})^\top \epsilon(\mathbf{w}^{(k)}) + 2(\mathbf{w} - \mathbf{w}^{(k)})^\top \mathbf{J}^\top \epsilon(\mathbf{w}^{(k)}) \\ &\quad + (\mathbf{w} - \mathbf{w}^{(k)})^\top \mathbf{J}^\top \mathbf{J} (\mathbf{w} - \mathbf{w}^{(k)}) \end{aligned}$$

# Trust Region

- Solution given by  $\nabla_{\mathbf{w}} L_{approx}(\mathbf{w}) = 0$  gives

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - (\mathbf{J}^T \mathbf{J})^{-1} \mathbf{J}^T \boldsymbol{\epsilon}(\mathbf{w}^{(k)})$$

- Can lead us in the wrong direction
- Instead use  $\mathbf{w}^{(k+1)} = \operatorname{argmin}_{\mathbf{w}} L_{approx}(\mathbf{w}) + \nu \|\mathbf{w} - \mathbf{w}^{(k)}\|^2$

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - (\mathbf{J}^T \mathbf{J} + \nu \mathbf{I})^{-1} \mathbf{J}^T \boldsymbol{\epsilon}(\mathbf{w})$$

- $\nu$  limits the step size
- If predicted reduction in error is accurate reduce  $\nu$ , if predicted reduction in error is very poor increase  $\nu$

# Trust Region

- Solution given by  $\nabla_{\mathbf{w}} L_{approx}(\mathbf{w}) = 0$  gives

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - (\mathbf{J}^T \mathbf{J})^{-1} \mathbf{J}^T \boldsymbol{\epsilon}(\mathbf{w}^{(k)})$$

- Can lead us in the wrong direction

- Instead use  $\mathbf{w}^{(k+1)} = \operatorname{argmin}_{\mathbf{w}} L_{approx}(\mathbf{w}) + \nu \|\mathbf{w} - \mathbf{w}^{(k)}\|^2$

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - (\mathbf{J}^T \mathbf{J} + \nu \mathbf{I})^{-1} \mathbf{J}^T \boldsymbol{\epsilon}(\mathbf{w})$$

- $\nu$  limits the step size
- If predicted reduction in error is accurate reduce  $\nu$ , if predicted reduction in error is very poor increase  $\nu$

# Trust Region

- Solution given by  $\nabla_{\mathbf{w}} L_{approx}(\mathbf{w}) = 0$  gives

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - (\mathbf{J}^T \mathbf{J})^{-1} \mathbf{J}^T \boldsymbol{\epsilon}(\mathbf{w}^{(k)})$$

- Can lead us in the wrong direction
- Instead use  $\mathbf{w}^{(k+1)} = \operatorname{argmin}_{\mathbf{w}} L_{approx}(\mathbf{w}) + \nu \|\mathbf{w} - \mathbf{w}^{(k)}\|^2$

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - (\mathbf{J}^T \mathbf{J} + \nu \mathbf{I})^{-1} \mathbf{J}^T \boldsymbol{\epsilon}(\mathbf{w})$$

- $\nu$  limits the step size
- If predicted reduction in error is accurate reduce  $\nu$ , if predicted reduction in error is very poor increase  $\nu$

# Trust Region

- Solution given by  $\nabla_{\mathbf{w}} L_{approx}(\mathbf{w}) = 0$  gives

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - (\mathbf{J}^T \mathbf{J})^{-1} \mathbf{J}^T \boldsymbol{\epsilon}(\mathbf{w}^{(k)})$$

- Can lead us in the wrong direction
- Instead use  $\mathbf{w}^{(k+1)} = \operatorname{argmin}_{\mathbf{w}} L_{approx}(\mathbf{w}) + \nu \|\mathbf{w} - \mathbf{w}^{(k)}\|^2$

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - (\mathbf{J}^T \mathbf{J} + \nu \mathbf{I})^{-1} \mathbf{J}^T \boldsymbol{\epsilon}(\mathbf{w})$$

- $\nu$  limits the step size
- If predicted reduction in error is accurate reduce  $\nu$ , if predicted reduction in error is very poor increase  $\nu$

# Trust Region

- Solution given by  $\nabla_{\mathbf{w}} L_{approx}(\mathbf{w}) = 0$  gives

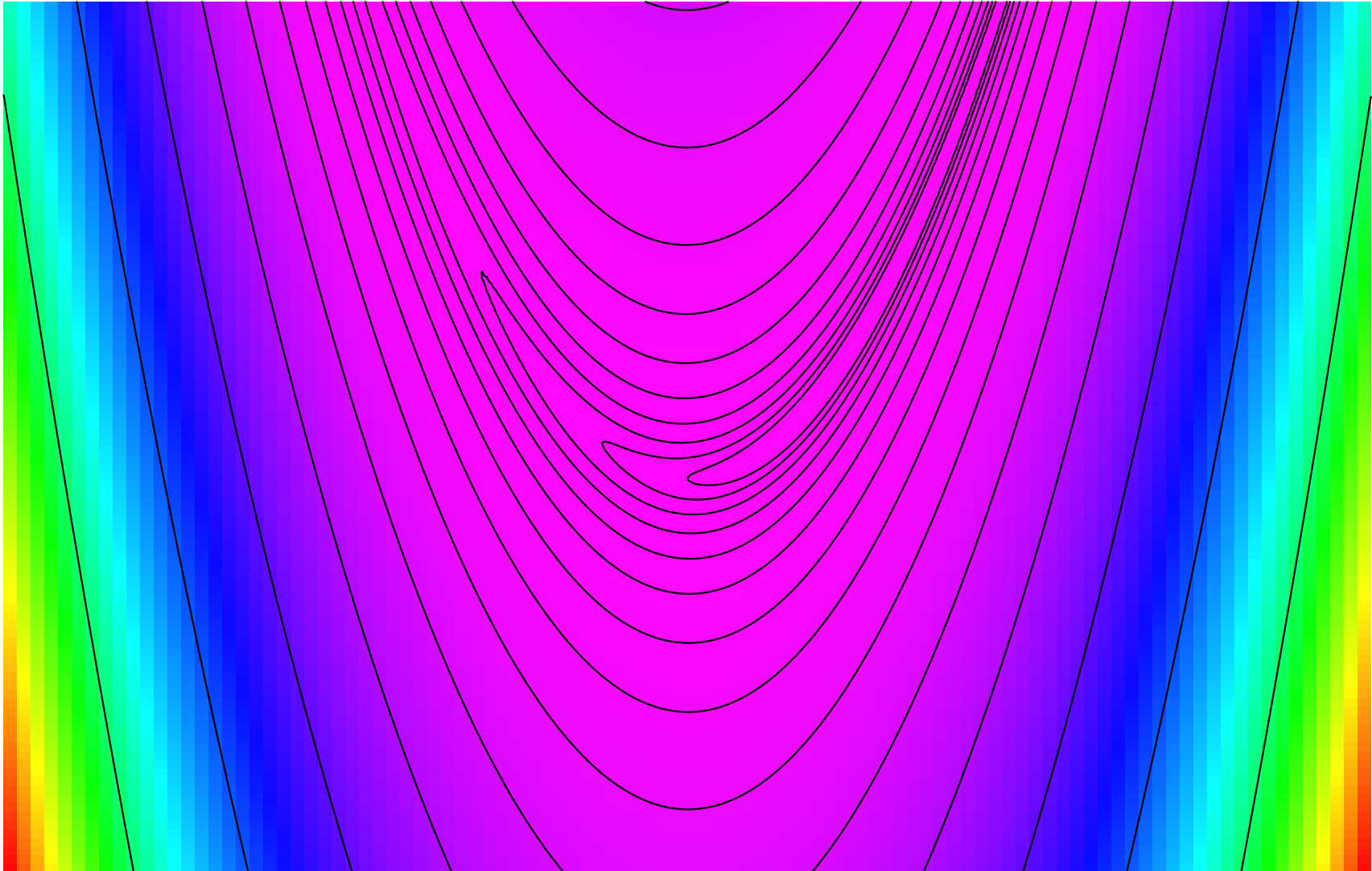
$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - (\mathbf{J}^T \mathbf{J})^{-1} \mathbf{J}^T \boldsymbol{\epsilon}(\mathbf{w}^{(k)})$$

- Can lead us in the wrong direction
- Instead use  $\mathbf{w}^{(k+1)} = \operatorname{argmin}_{\mathbf{w}} L_{approx}(\mathbf{w}) + \nu \|\mathbf{w} - \mathbf{w}^{(k)}\|^2$

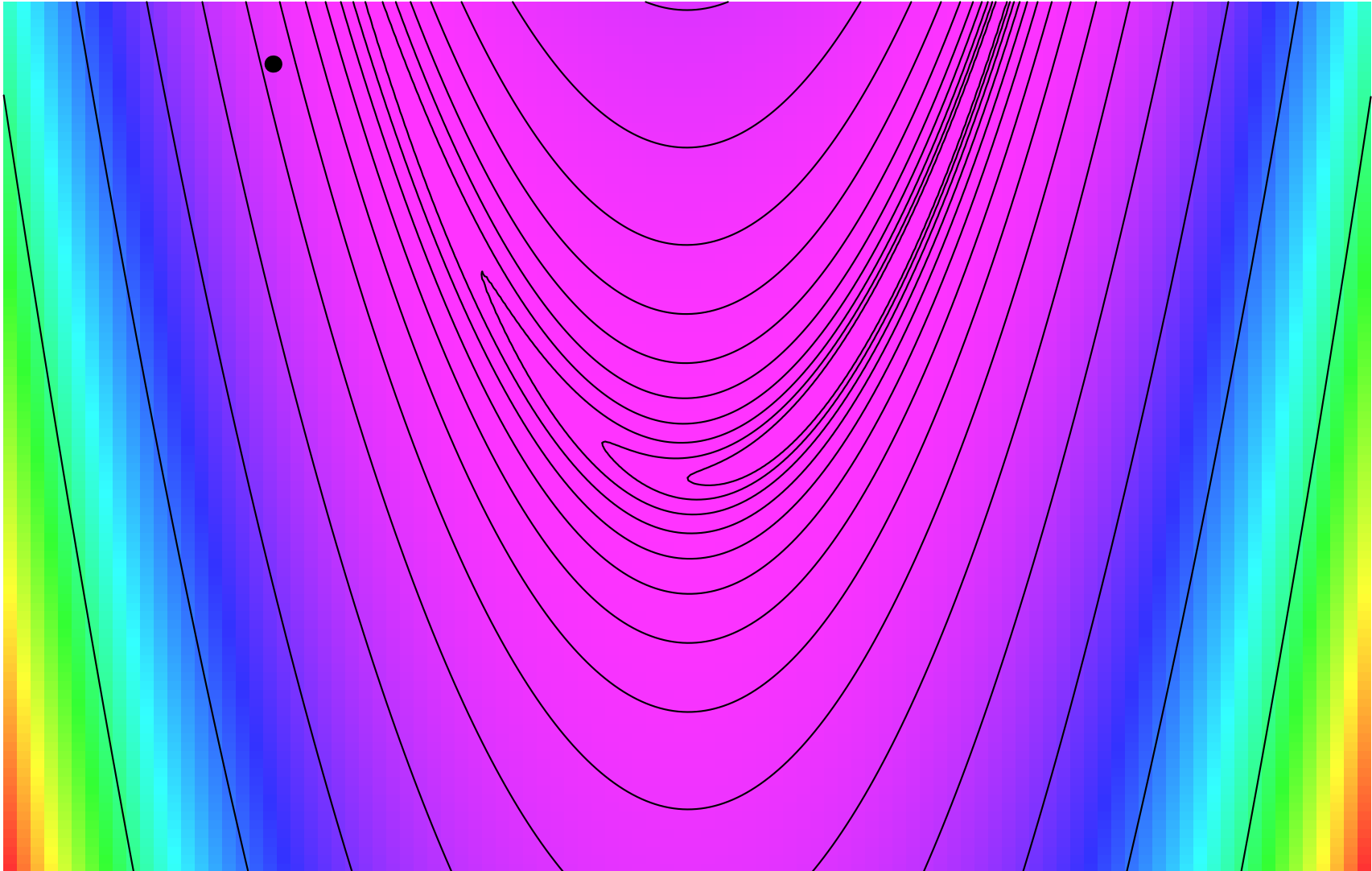
$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - (\mathbf{J}^T \mathbf{J} + \nu \mathbf{I})^{-1} \mathbf{J}^T \boldsymbol{\epsilon}(\mathbf{w})$$

- $\nu$  limits the step size
- If predicted reduction in error is accurate reduce  $\nu$ , if predicted reduction in error is very poor increase  $\nu$

$$\epsilon_1 = 10(x_2 - x_1^2) \text{ and } \epsilon_2 = 1 - x_1$$

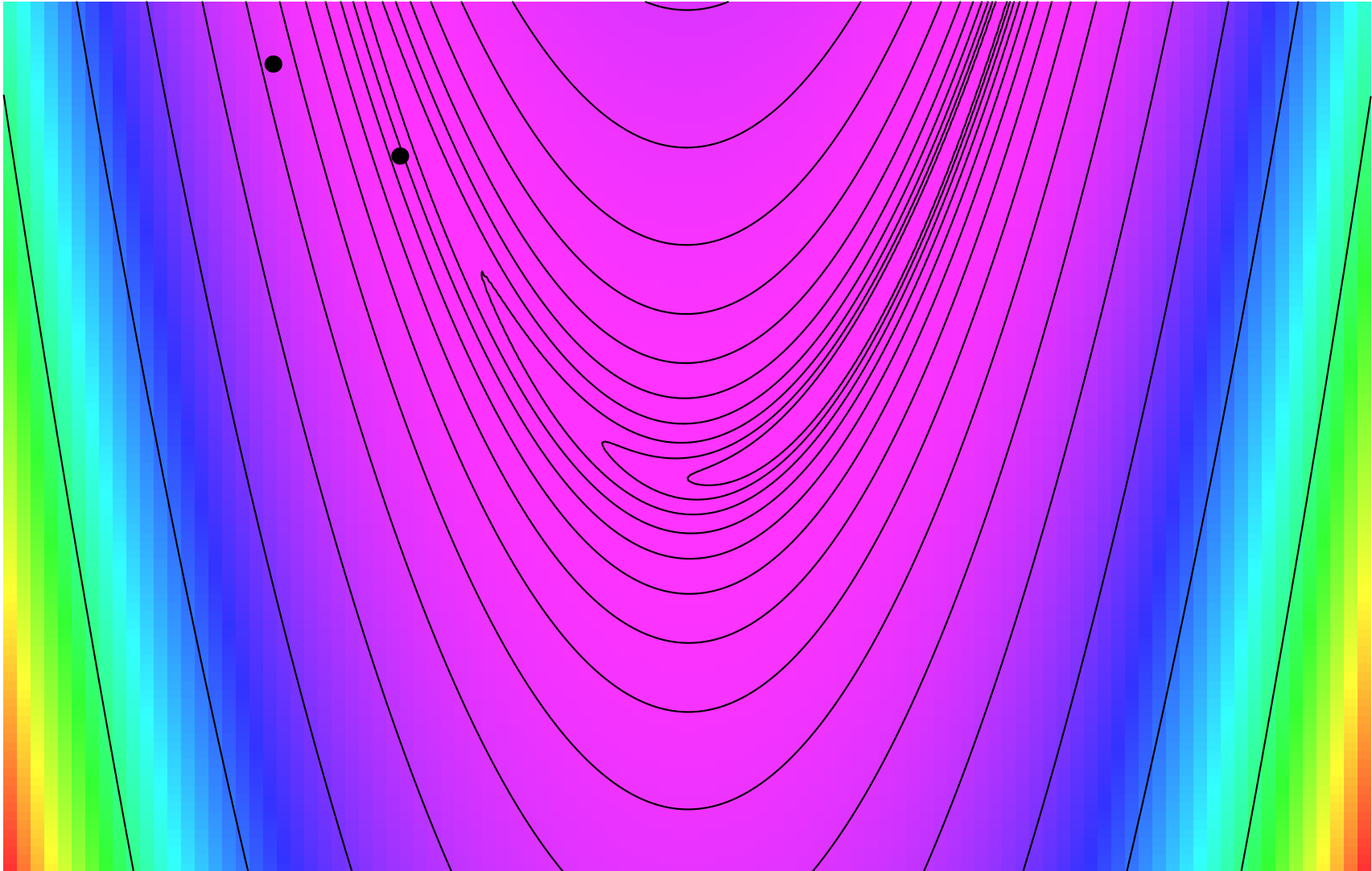


$$\epsilon_1 = 10(x_2 - x_1^2) \text{ and } \epsilon_2 = 1 - x_1$$

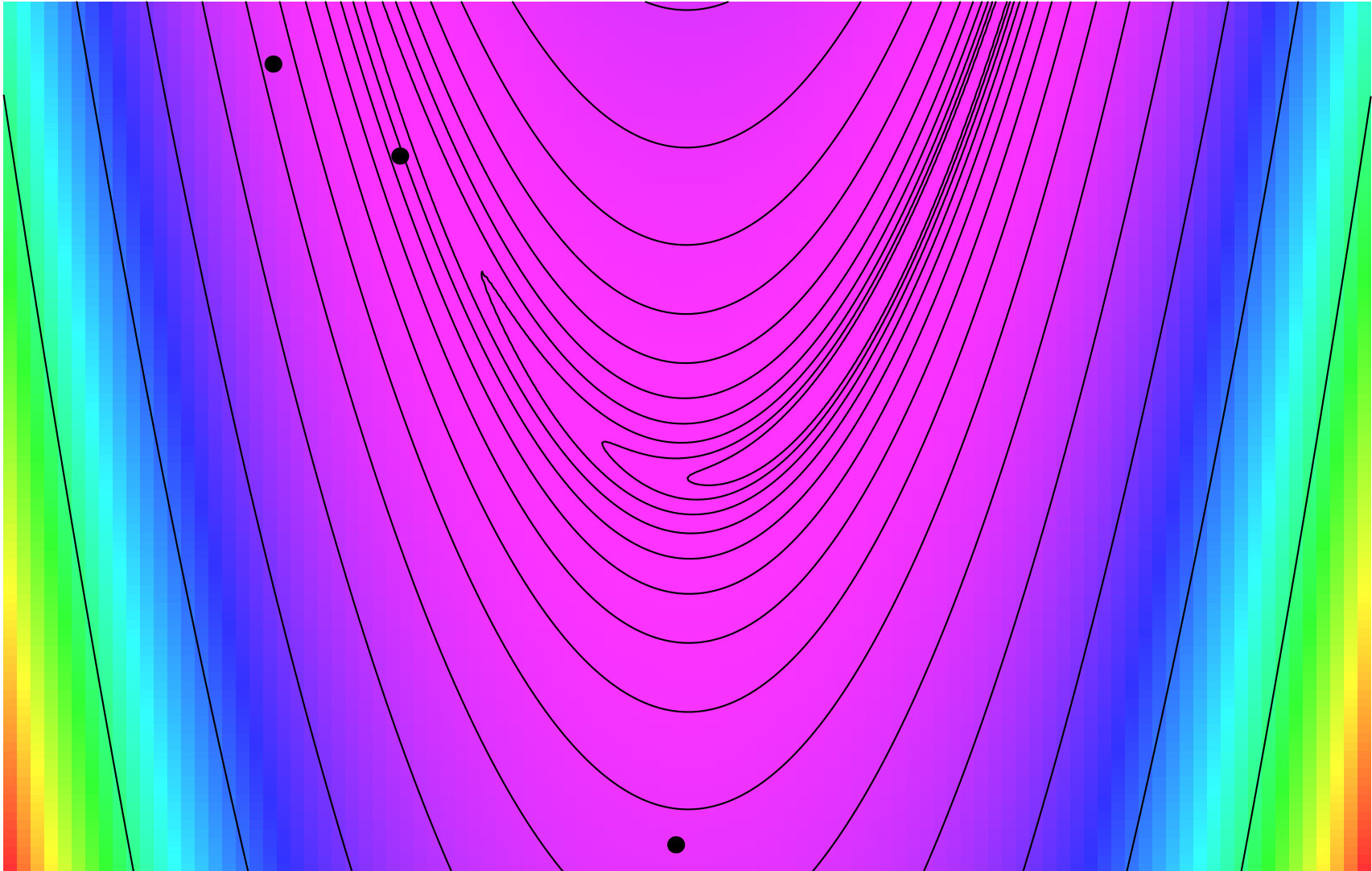




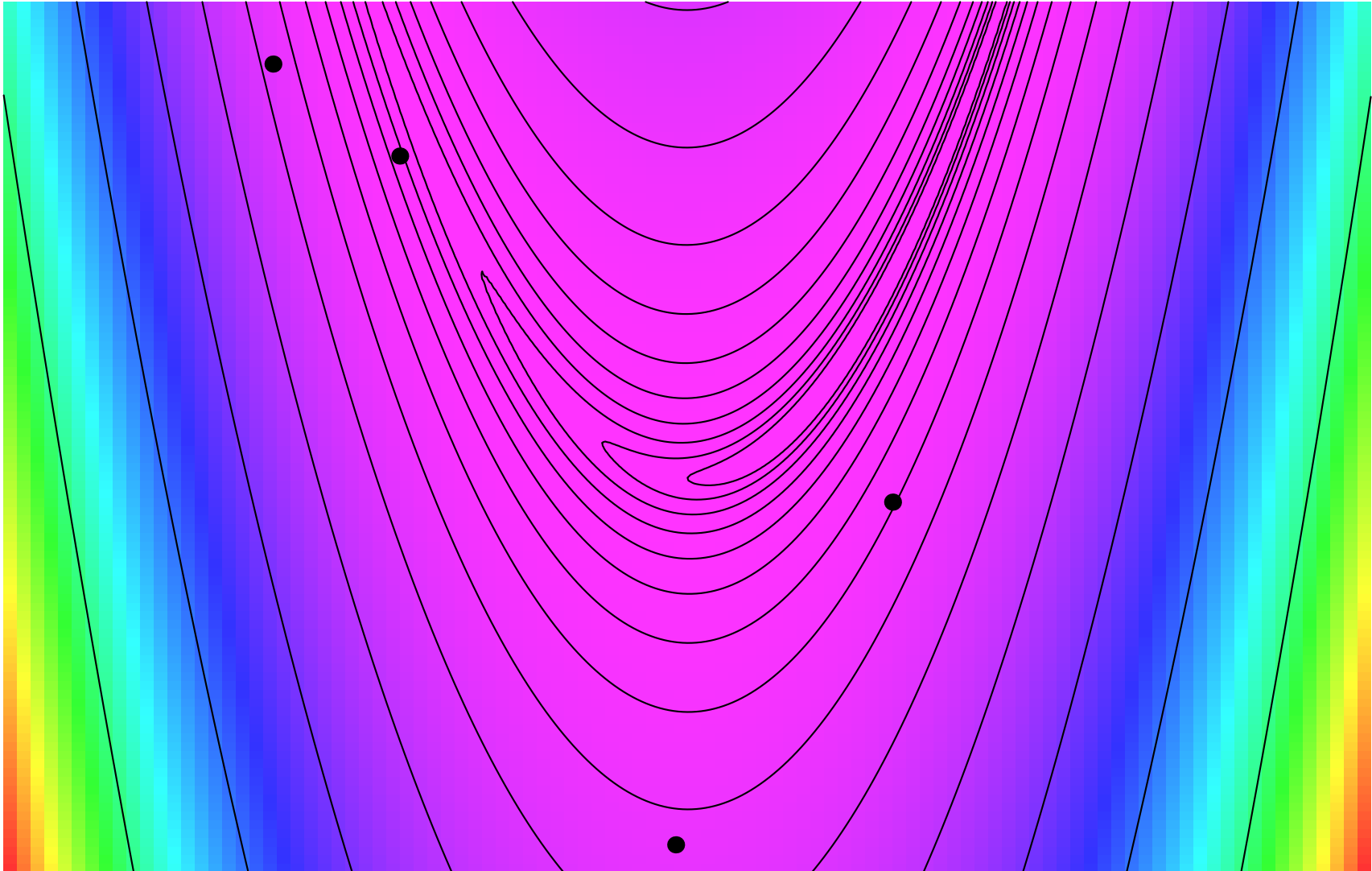
$$\epsilon_1 = 10(x_2 - x_1^2) \text{ and } \epsilon_2 = 1 - x_1$$



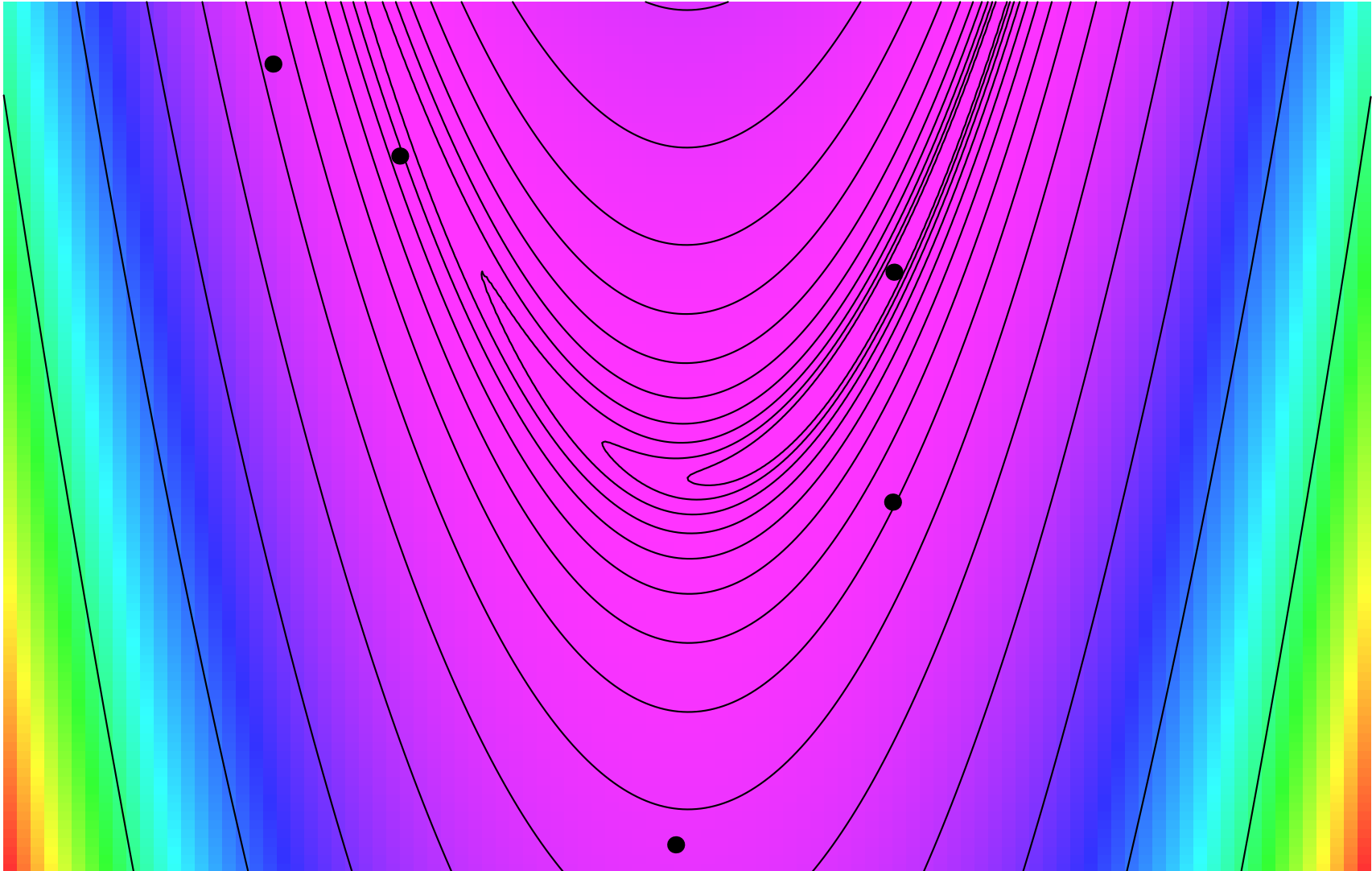
$$\epsilon_1 = 10(x_2 - x_1^2) \text{ and } \epsilon_2 = 1 - x_1$$



$$\epsilon_1 = 10(x_2 - x_1^2) \text{ and } \epsilon_2 = 1 - x_1$$

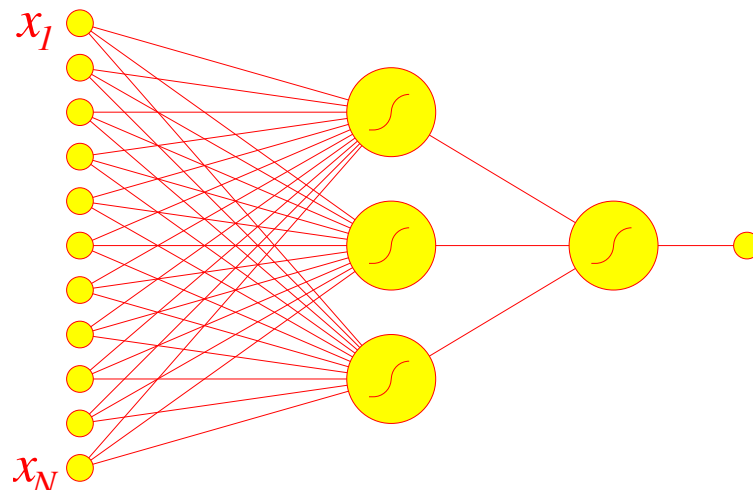


$$\epsilon_1 = 10(x_2 - x_1^2) \text{ and } \epsilon_2 = 1 - x_1$$



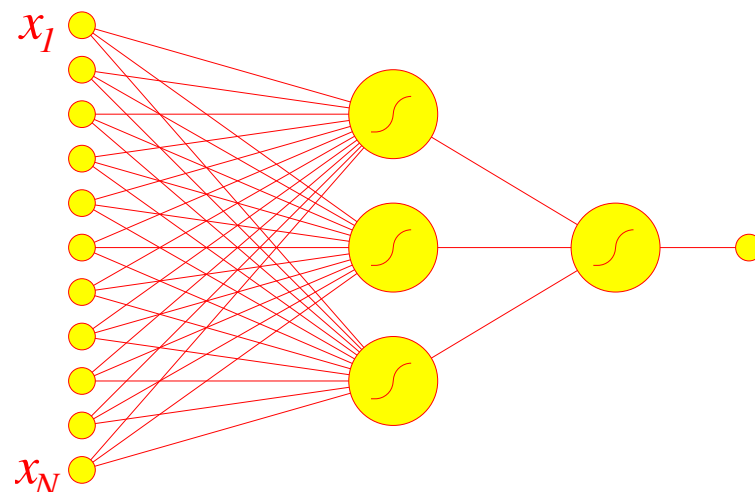
# Landscape Problems

- For sigmoid response function such as the logistic or tanh function there are very flat regions away from the bias
- Very little gradient information (can use momentum)
- Many local minima—not guaranteed to find global minimum
- Weight space is symmetric because of permutation symmetry



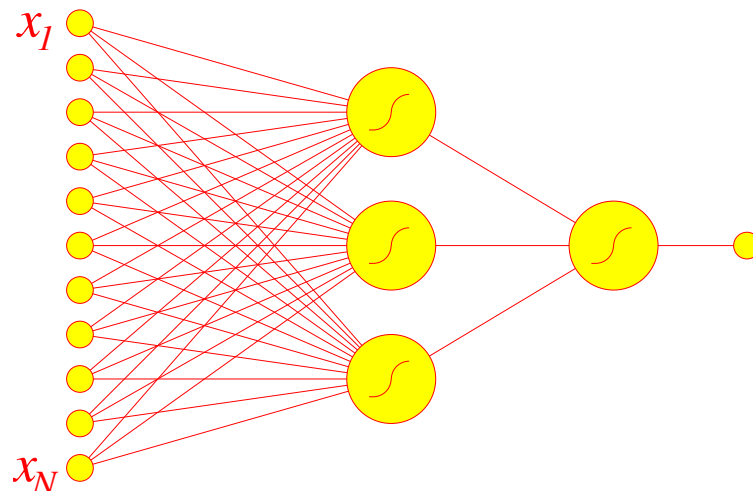
# Landscape Problems

- For sigmoid response function such as the logistic or tanh function there are very flat regions away from the bias
- Very little gradient information (can use momentum)
- Many local minima—not guaranteed to find global minimum
- Weight space is symmetric because of permutation symmetry



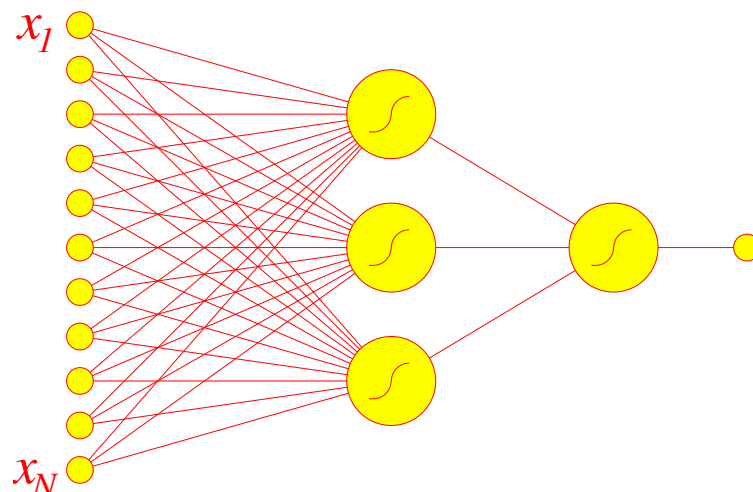
# Landscape Problems

- For sigmoid response function such as the logistic or tanh function there are very flat regions away from the bias
- Very little gradient information (can use momentum)
- **Many local minima**—not guaranteed to find global minimum
- Weight space is symmetric because of permutation symmetry



# Landscape Problems

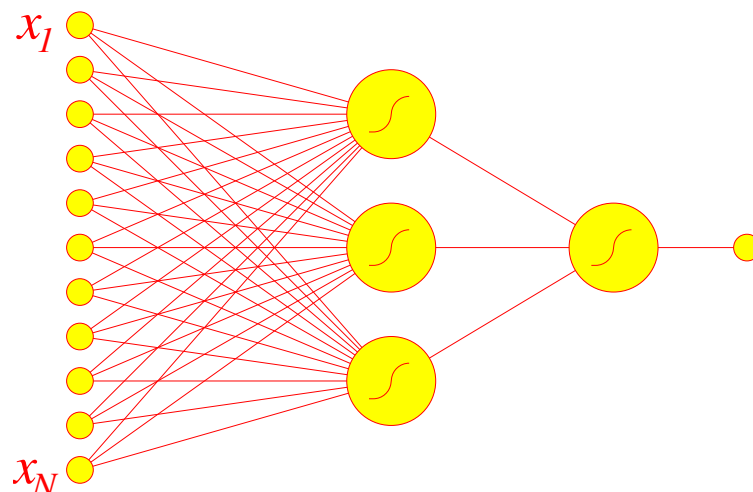
- For sigmoid response function such as the logistic or tanh function there are very flat regions away from the bias
- Very little gradient information (can use momentum)
- Many local minima—not guaranteed to find global minimum
- Weight space is symmetric because of permutation symmetry





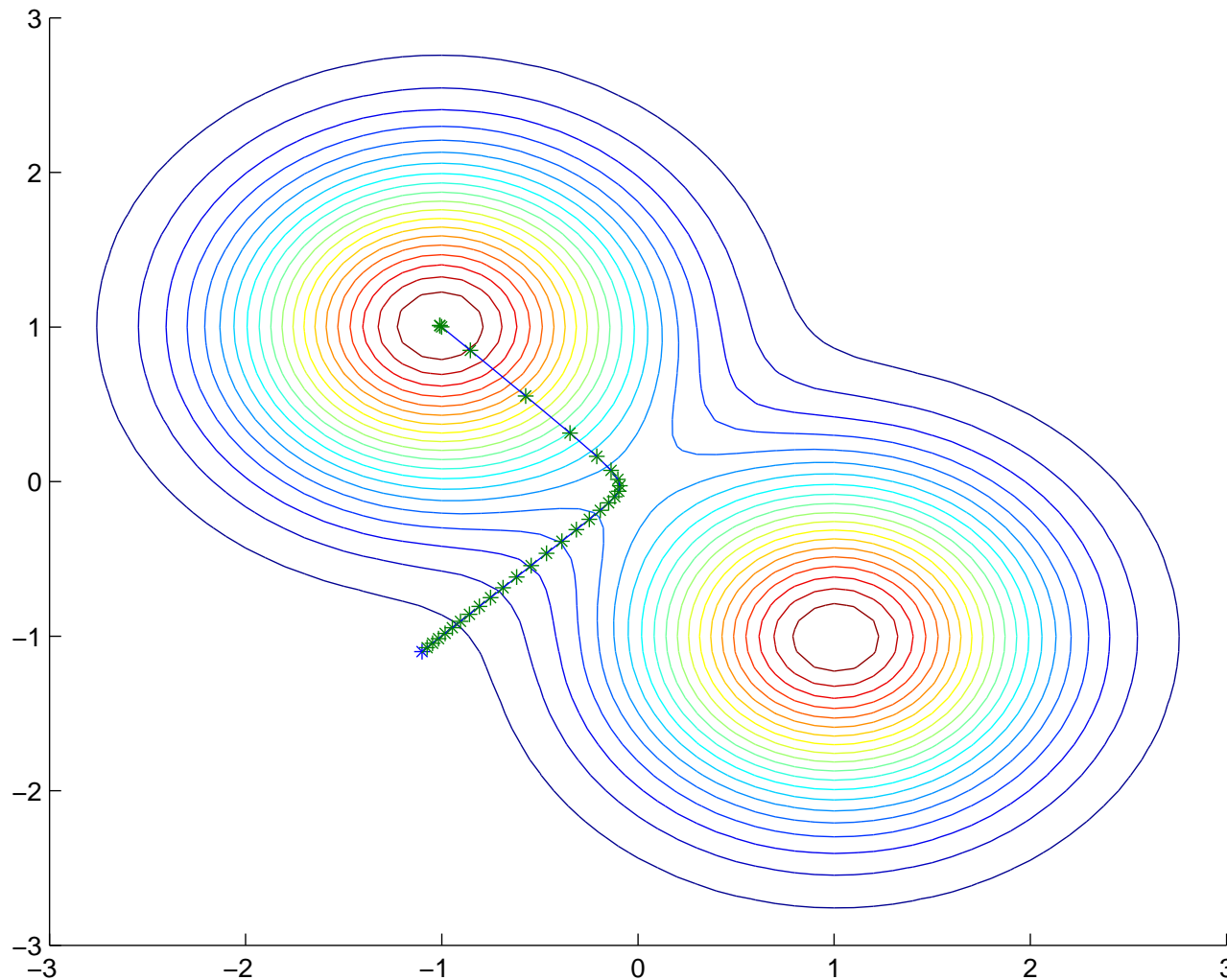
# Landscape Problems

- For sigmoid response function such as the logistic or tanh function there are very flat regions away from the bias
- Very little gradient information (can use momentum)
- Many local minima—not guaranteed to find global minimum
- Weight space is symmetric because of permutation symmetry



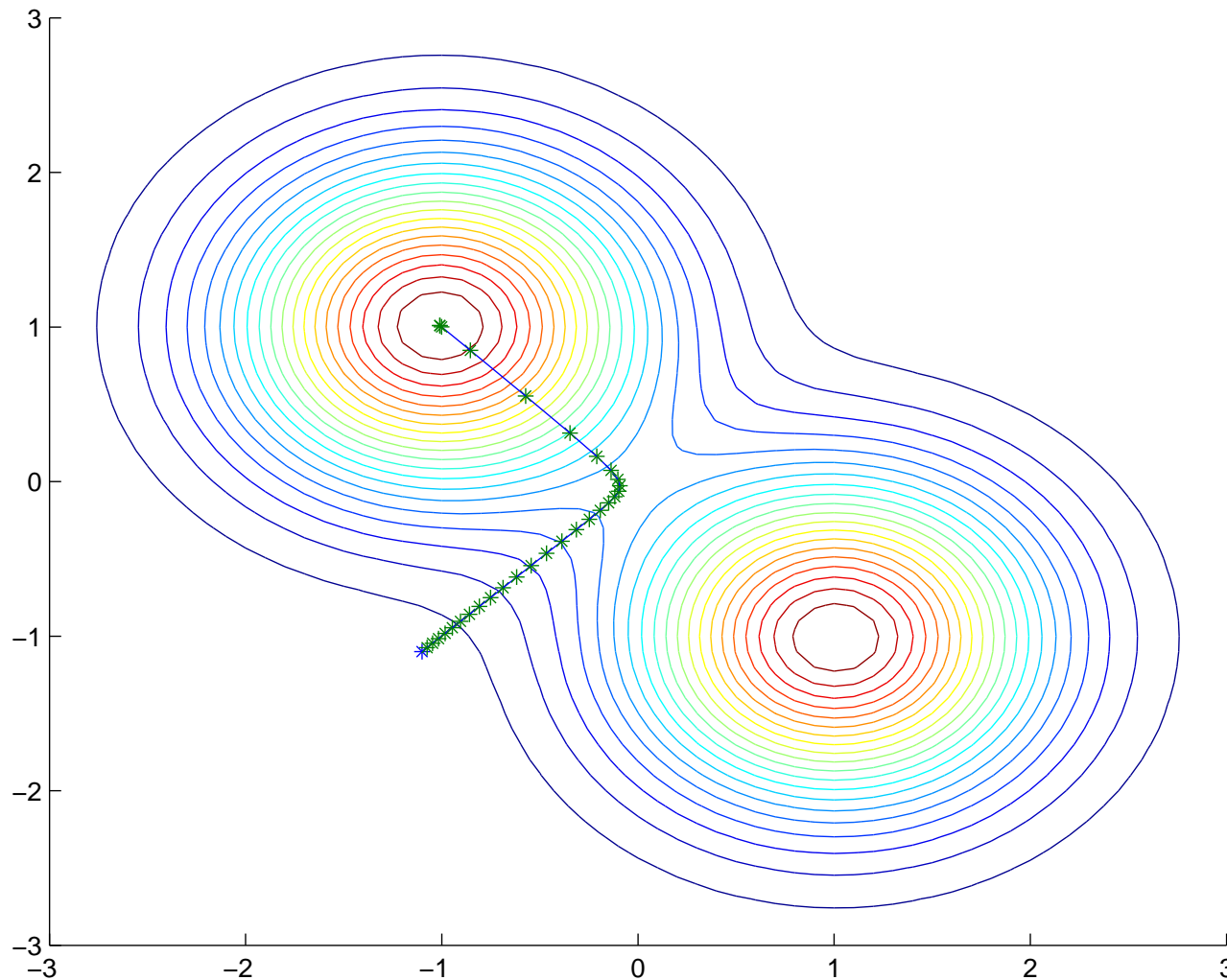
# Saddle-Points

- Many saddle-points due to permutation symmetry of hidden nodes



# Saddle-Points

- Many saddle-points due to permutation symmetry of hidden nodes



# Modern Machine Learning

- Modern deep learning machines usually have so many parameters that it would be impractical to even estimate the Hessian
- We now commonly use ReLU's rather than sigmoids so the cost landscape is non-analytic
- Making assumptions about the minima being quadratic is no longer true
- In deep learning it is usually important to ensure that loss is differentiable almost everywhere

# Modern Machine Learning

- Modern deep learning machines usually have so many parameters that it would be impractical to even estimate the Hessian
- We now commonly use ReLU's rather than sigmoids so the cost landscape is non-analytic
- Making assumptions about the minima being quadratic is no longer true
- In deep learning it is usually important to ensure that loss is differentiable almost everywhere

# Modern Machine Learning

- Modern deep learning machines usually have so many parameters that it would be impractical to even estimate the Hessian
- We now commonly use ReLU's rather than sigmoids so the cost landscape is non-analytic
- Making assumptions about the minima being quadratic is no longer true
- In deep learning it is usually important to ensure that loss is differentiable almost everywhere

# Modern Machine Learning

- Modern deep learning machines usually have so many parameters that it would be impractical to even estimate the Hessian
- We now commonly use ReLU's rather than sigmoids so the cost landscape is non-analytic
- Making assumptions about the minima being quadratic is no longer true
- In deep learning it is usually important to ensure that loss is differentiable almost everywhere

# Modern Machine Learning

- Modern deep learning machines usually have so many parameters that it would be impractical to even estimate the Hessian
- We now commonly use ReLU's rather than sigmoids so the cost landscape is non-analytic
- Making assumptions about the minima being quadratic is no longer true
- In deep learning it is usually important to ensure that loss is differentiable almost everywhere **otherwise we can't use gradient descent**



# Summary

- Minimising the error gives a general purpose learning algorithm
  - ★ E.g. Logistic Perceptron, Multi-Layer Perceptron, CNNs, . . .
- The minimisation is in a high dimensional space
- Minimisation can be very time consuming
- Nasty things can happen, but there are many standard algorithms which work well
- In complex models there may be many minima—you are not guaranteed to find the best one

# Summary

- Minimising the error gives a general purpose learning algorithm
  - ★ E.g. Logistic Perceptron, Multi-Layer Perceptron, CNNs, . . .
- The minimisation is in a high dimensional space
- Minimisation can be very time consuming
- Nasty things can happen, but there are many standard algorithms which work well
- In complex models there may be many minima—you are not guaranteed to find the best one

# Summary

- Minimising the error gives a general purpose learning algorithm
  - ★ E.g. Logistic Perceptron, Multi-Layer Perceptron, CNNs, . . .
- The minimisation is in a high dimensional space
- Minimisation can be very time consuming
- Nasty things can happen, but there are many standard algorithms which work well
- In complex models there may be many minima—you are not guaranteed to find the best one

# Summary

- Minimising the error gives a general purpose learning algorithm
  - ★ E.g. Logistic Perceptron, Multi-Layer Perceptron, CNNs, . . .
- The minimisation is in a high dimensional space
- Minimisation can be very time consuming
- Nasty things can happen, but there are many standard algorithms which work well
- In complex models there may be many minima—you are not guaranteed to find the best one

# Summary

- Minimising the error gives a general purpose learning algorithm
  - ★ E.g. Logistic Perceptron, Multi-Layer Perceptron, CNNs, . . .
- The minimisation is in a high dimensional space
- Minimisation can be very time consuming
- Nasty things can happen, but there are many standard algorithms which work well
- In complex models there may be many minima—you are not guaranteed to find the best one

# Summary

- Minimising the error gives a general purpose learning algorithm
  - ★ E.g. Logistic Perceptron, Multi-Layer Perceptron, CNNs, . . .
- The minimisation is in a high dimensional space
- Minimisation can be very time consuming
- Nasty things can happen, but there are many standard algorithms which work well
- In complex models there may be many minima—you are not guaranteed to find the best one