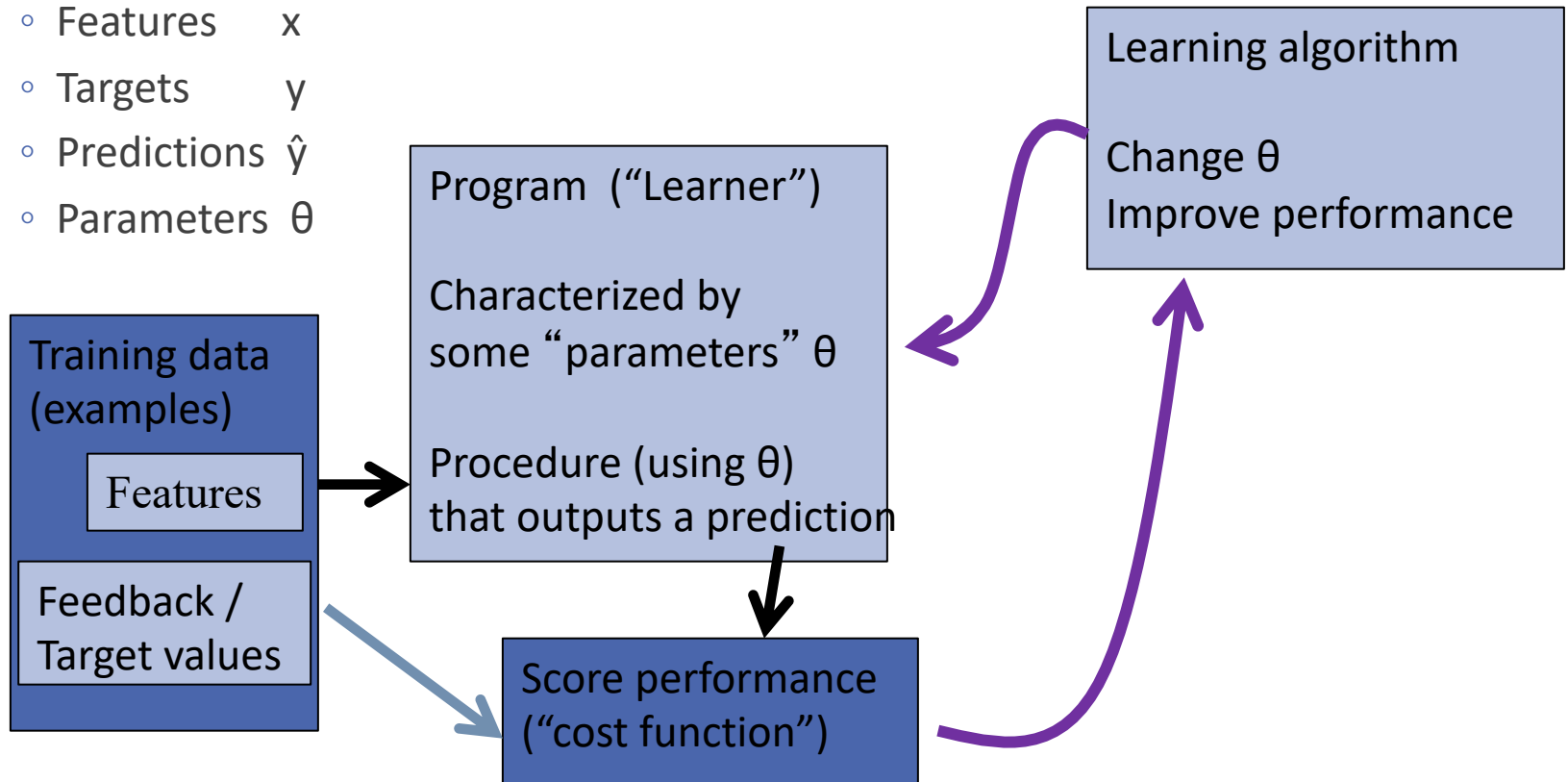


# Supervised learning

## Notation

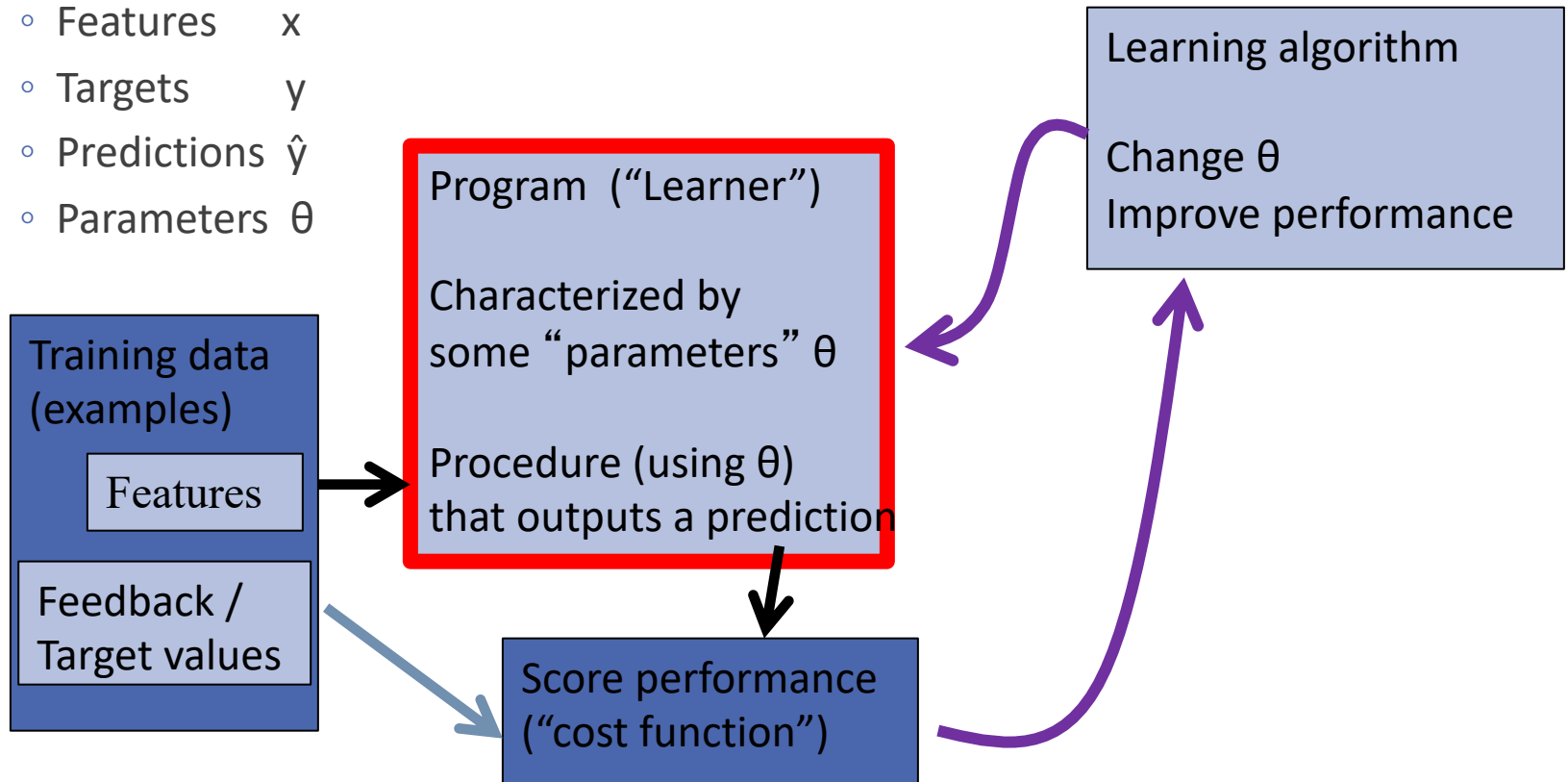
- Features  $x$
- Targets  $y$
- Predictions  $\hat{y}$
- Parameters  $\theta$



# Supervised learning

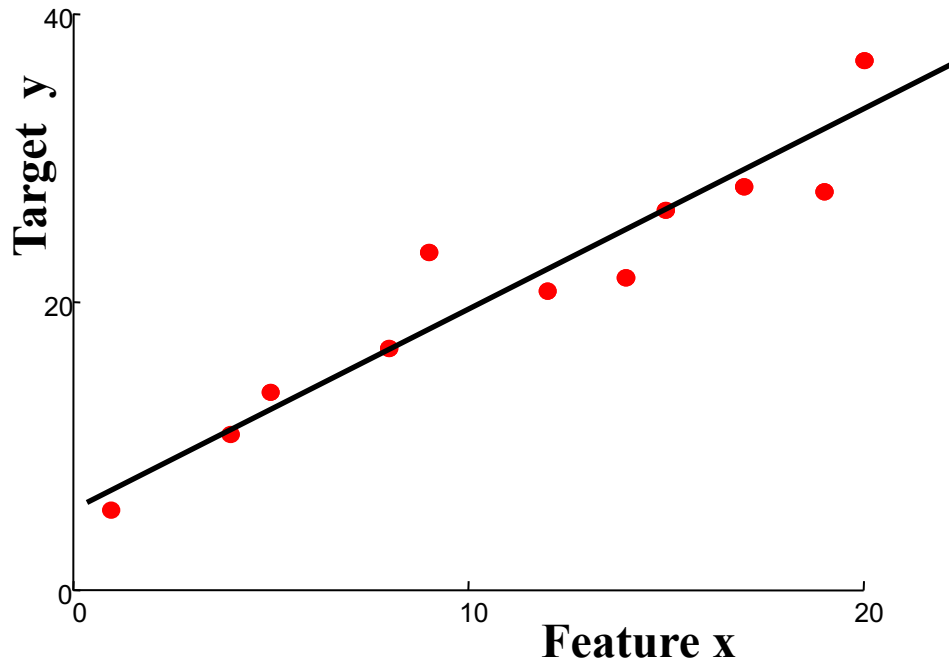
## Notation

- Features  $x$
- Targets  $y$
- Predictions  $\hat{y}$
- Parameters  $\theta$



# Linear regression

---



**“Predictor”:**

Evaluate line:

$$r = \theta_0 + \theta_1 x_1$$

return r

Define form of function  $f(x)$  explicitly

Find a good  $f(x)$  within that family

# Notation

---

$$\hat{y}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots$$

Define “feature”  $x_0 = 1$  (constant)

Then

$$\hat{y}(x) = \theta x^T$$

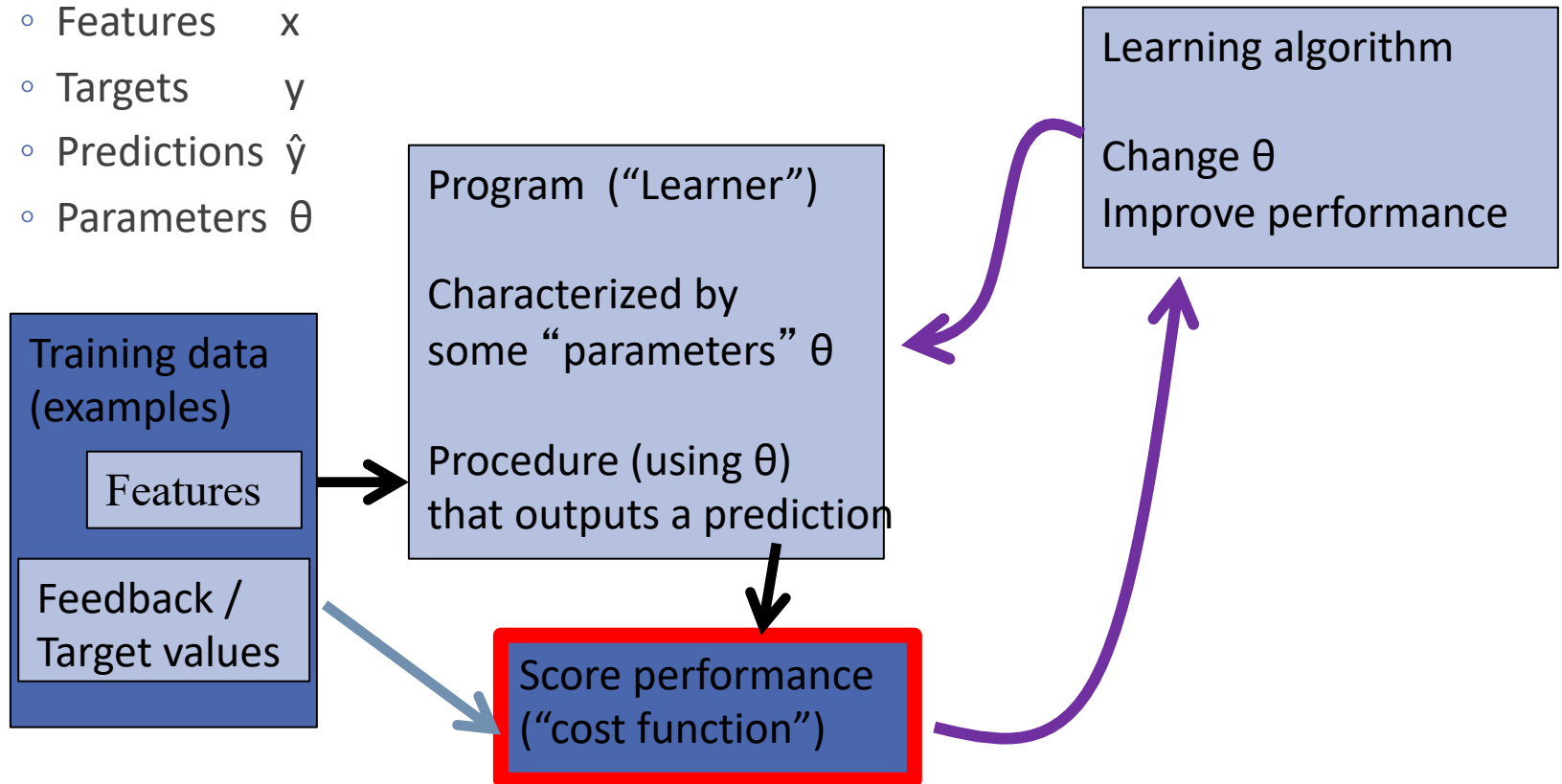
$$\underline{\theta} = [\theta_0, \dots, \theta_n]$$

$$\underline{x} = [1, x_1, \dots, x_n]$$

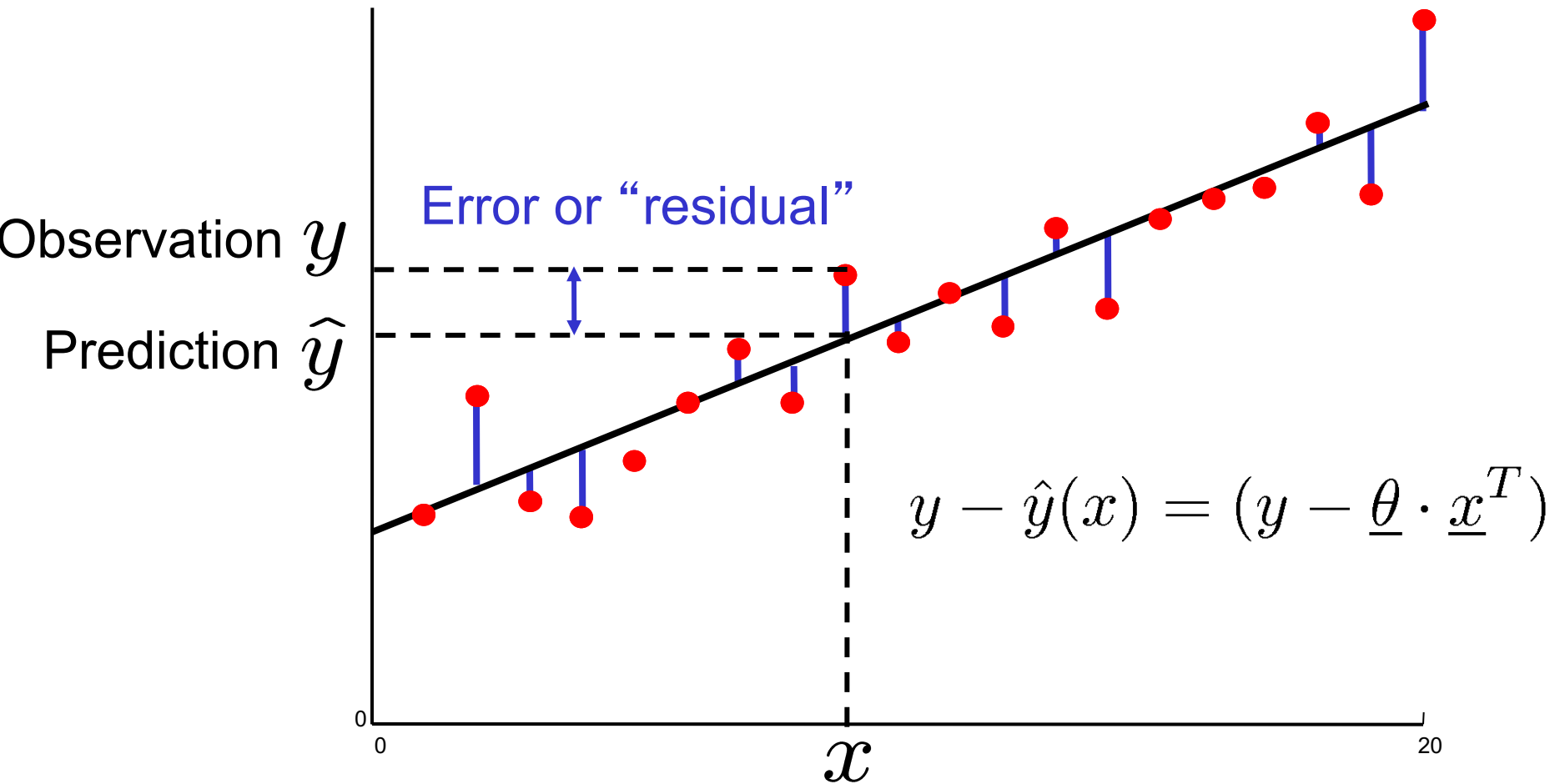
# Supervised learning

## Notation

- Features  $x$
- Targets  $y$
- Predictions  $\hat{y}$
- Parameters  $\theta$



# Measuring error



# Mean squared error

---

How can we quantify the error?

$$\begin{aligned}\text{MSE, } J(\underline{\theta}) &= \frac{1}{m} \sum_j (y^{(j)} - \hat{y}(x^{(j)}))^2 \\ &= \frac{1}{m} \sum_j (y^{(j)} - \underline{\theta} \cdot \underline{x}^{(j)T})^2\end{aligned}$$

Why choosing exactly this error measure and not something else?

- Computationally convenient (more later)
- Measures the variance of the residuals
- Corresponds to likelihood under Gaussian model of “noise”

$$N(y; \hat{y}, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2} (y - \hat{y})^2\right\}$$

# MSE Cost function

---

$$\begin{aligned}\text{MSE, } J(\underline{\theta}) &= \frac{1}{m} \sum_j (y^{(j)} - \hat{y}(x^{(j)}))^2 \\ &= \frac{1}{m} \sum_j (y^{(j)} - \underline{\theta} \cdot \underline{x}^{(j)T})^2\end{aligned}$$

$$\underline{\theta} = [\theta_0, \dots, \theta_n]$$

$$\underline{y} = \begin{bmatrix} y^{(1)} & \dots & y^{(m)} \end{bmatrix}^T$$

$$\underline{X} = \begin{bmatrix} x_0^{(1)} & \dots & x_n^{(1)} \\ \vdots & \ddots & \vdots \\ x_0^{(m)} & \dots & x_n^{(m)} \end{bmatrix}$$

$$J(\underline{\theta}) = \frac{1}{m} (\underline{y}^T - \underline{\theta} \underline{X}^T) \cdot (\underline{y}^T - \underline{\theta} \underline{X}^T)^T$$

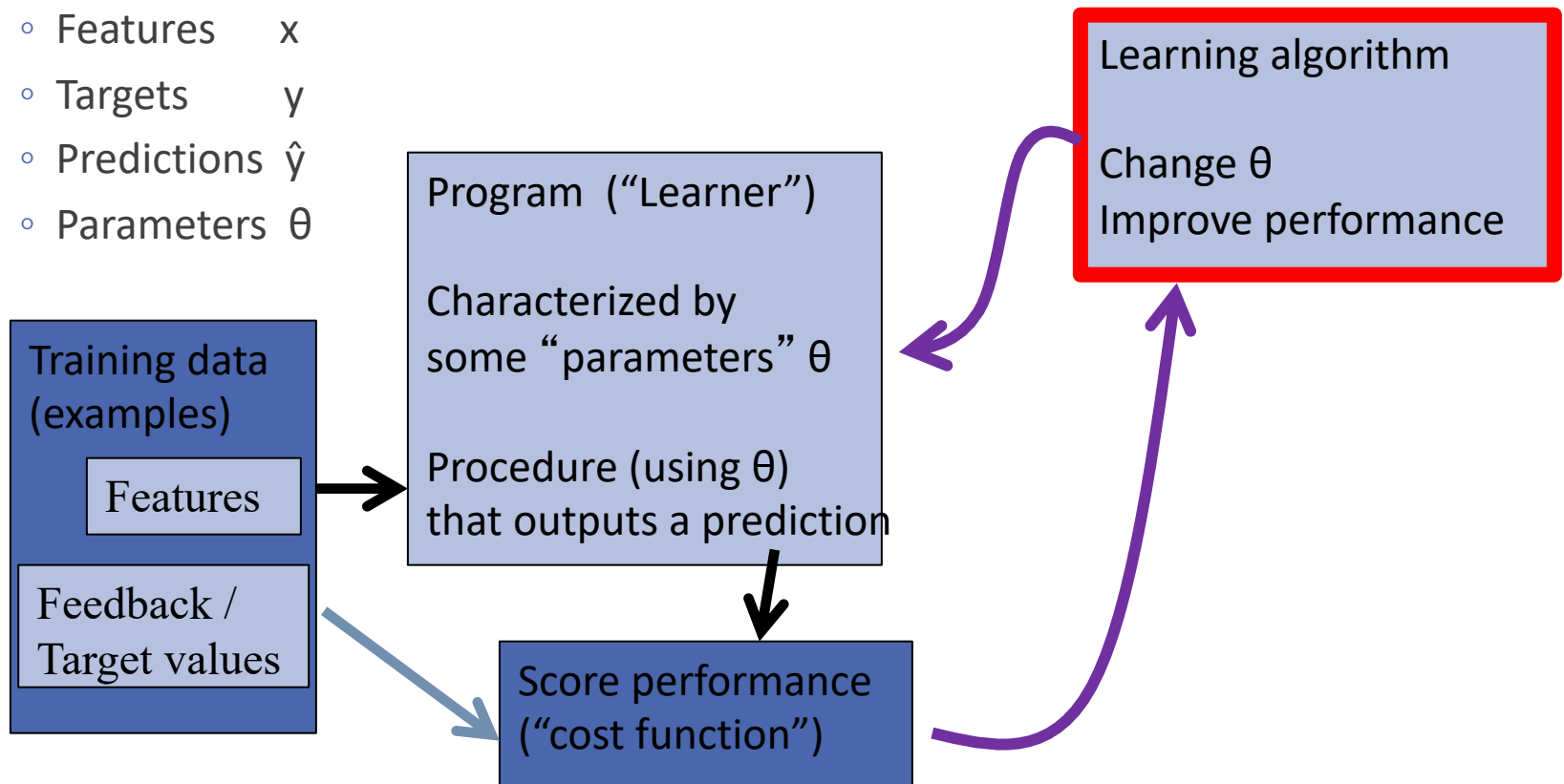
```
# Python / NumPy:  
e = Y - X.dot( theta.T );  
J = e.T.dot( e ) / m # = np.mean( e ** 2 )
```



# Supervised learning

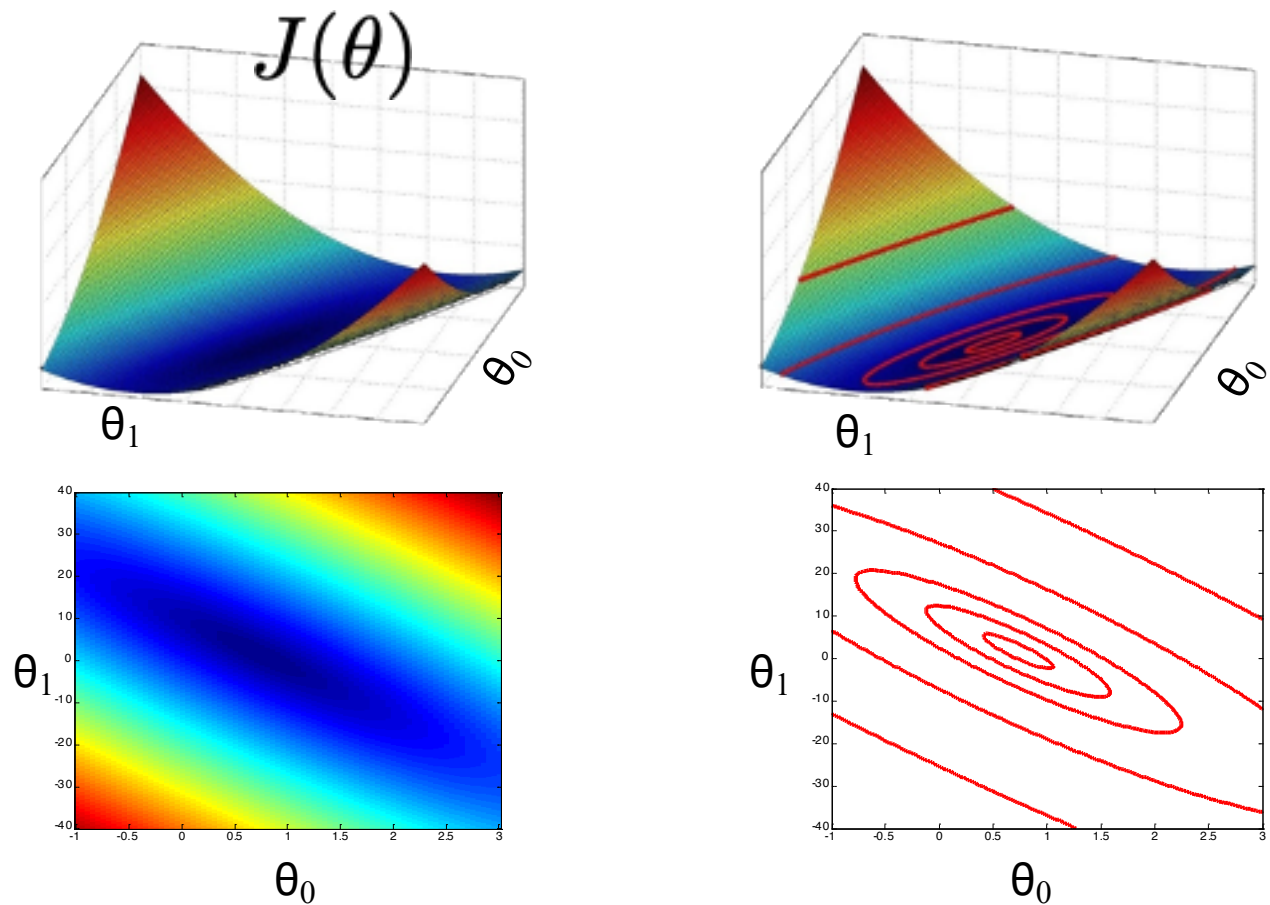
## Notation

- Features  $x$
- Targets  $y$
- Predictions  $\hat{y}$
- Parameters  $\theta$



# Visualizing the cost function

---

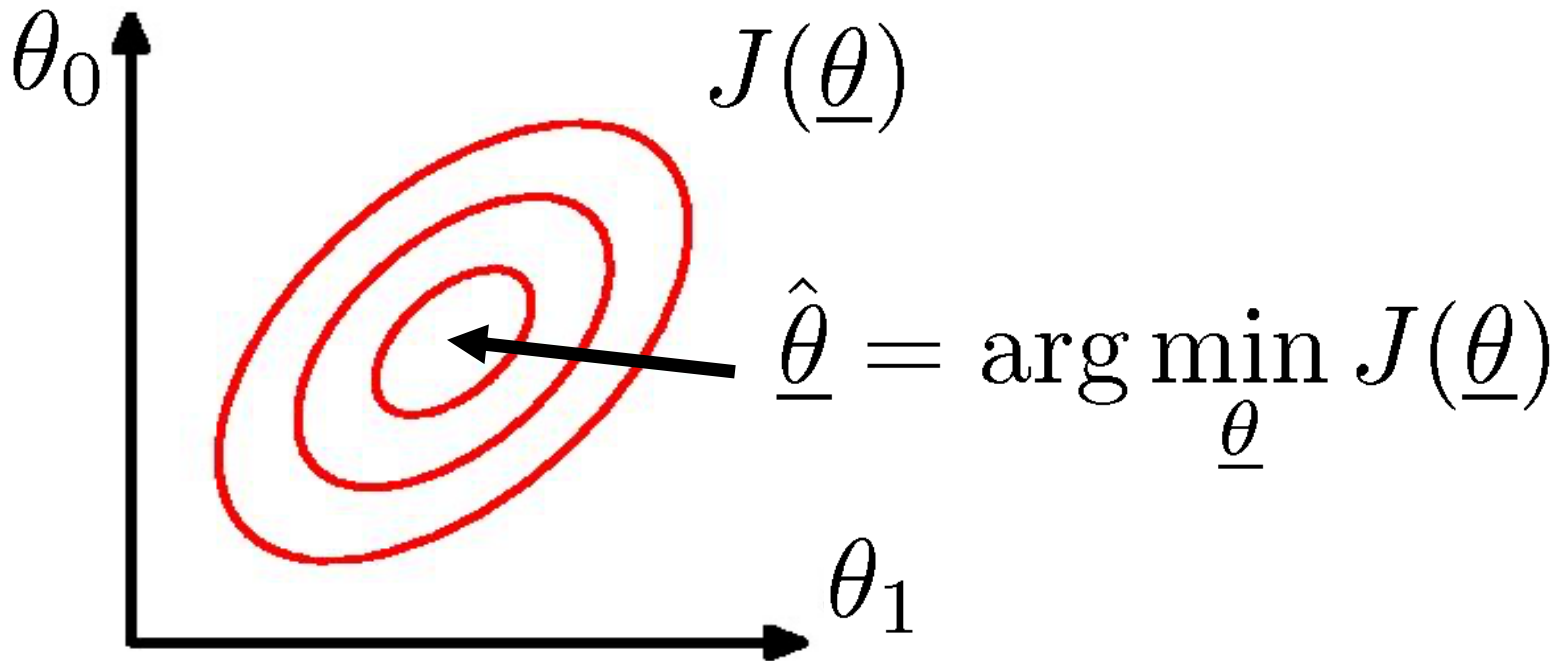


# Finding good parameters

---

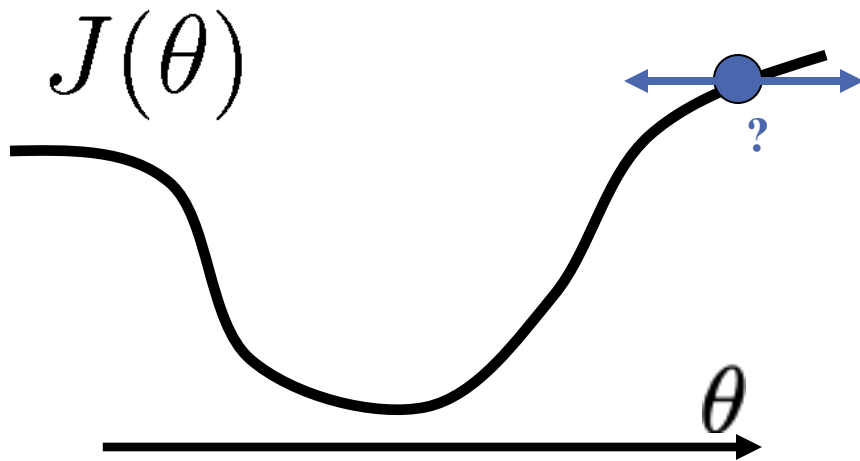
Want to find parameters which minimize our error...

Think of a cost “surface”: error residual for that  $\theta$ ...



# Gradient descent

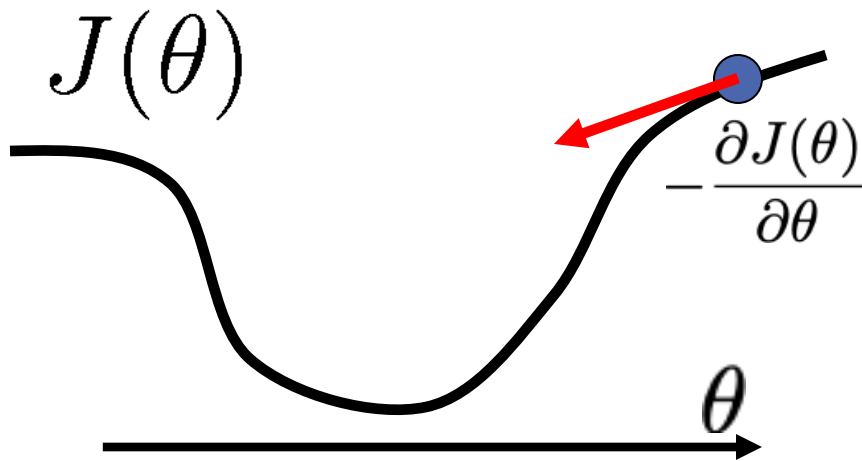
---



- How to change  $\theta$  to improve  $J(\theta)$ ?
- Choose a direction in which  $J(\theta)$  is decreasing

# Gradient descent

---



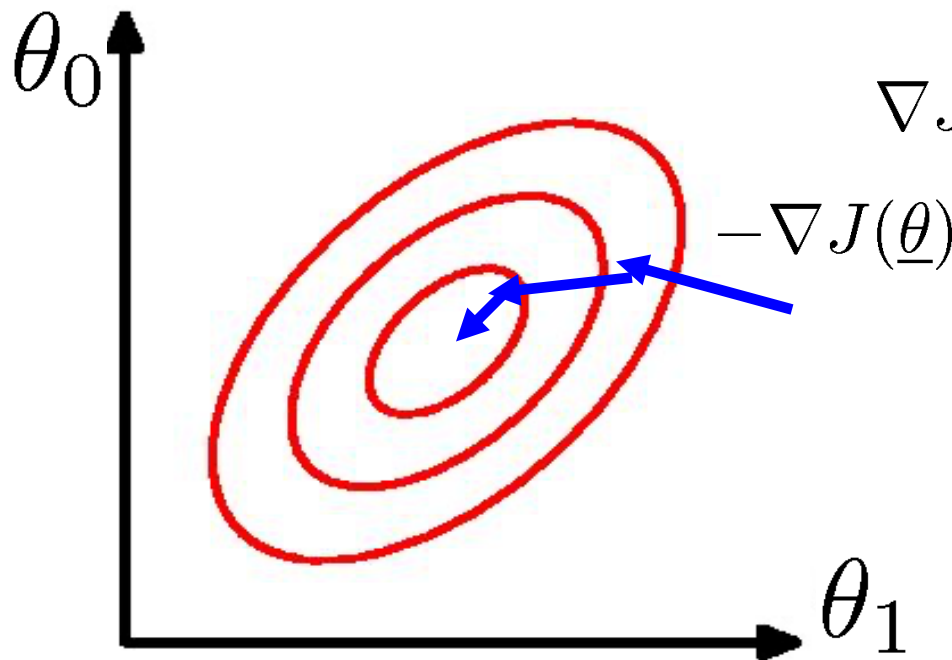
- How to change  $\theta$  to improve  $J(\theta)$ ?
- Choose a direction in which  $J(\theta)$  is decreasing
- Derivative  $\frac{\partial J(\theta)}{\partial \theta}$
- Positive  $\Rightarrow$  increasing
- Negative  $\Rightarrow$  decreasing

# Gradient descent in >2 dimensions

---

- Gradient vector

$$\nabla J(\underline{\theta}) = \left[ \frac{\partial J(\underline{\theta})}{\partial \theta_0} \quad \frac{\partial J(\underline{\theta})}{\partial \theta_1} \quad \dots \right]$$



Indicates direction of steepest ascent  
(negative = steepest descent)

# Gradient descent

---

Initialization

Initialize  $\theta$

Step size  $\alpha$

Do{

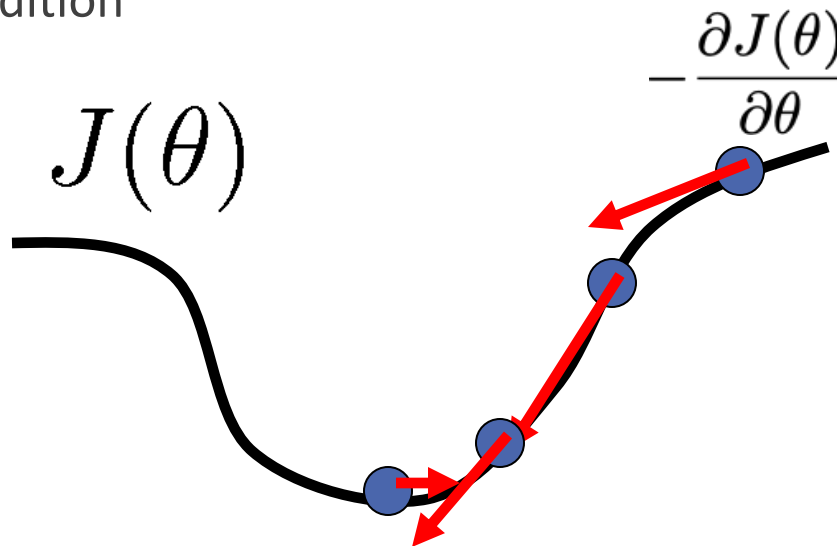
- Can change as a function of iteration

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} J(\theta)$$

Gradient direction

} while ( $\alpha \|\nabla_{\theta} J\| > \epsilon$ )

Stopping condition



# Gradient for the MSE

MSE

$$J(\underline{\theta}) = \frac{1}{m} \sum_j (y^{(j)} - \underline{\theta} \cdot \underline{x}^{(j)T})^2$$

$\nabla J = ?$

$$J(\underline{\theta}) = \frac{1}{m} \sum_j \overbrace{(y^{(j)} - \theta_0 \underline{x}_0^{(j)} - \theta_1 \underline{x}_1^{(j)} - \dots)}^{e_j(\theta)}^2$$

$$\frac{\partial J}{\partial \theta_0} = \frac{\partial}{\partial \theta_0} \frac{1}{m} \sum_j (e_j(\theta))^2$$

$$= \frac{1}{m} \sum_j \frac{\partial}{\partial \theta_0} (e_j(\theta))^2$$

$$= \frac{1}{m} \sum_j 2e_j(\theta) \frac{\partial}{\partial \theta_0} e_j(\theta)$$

$$\frac{\partial}{\partial \theta_0} e_j(\theta) = \cancel{\frac{\partial}{\partial \theta_0} y^{(j)}} - \frac{\partial}{\partial \theta_0} \theta_0 x_0^{(j)} - \cancel{\frac{\partial}{\partial \theta_0} \theta_1 x_1^{(j)}} - \dots$$

**0** **0**

$$= -x_0^{(j)}$$



# Gradient for the MSE

---

MSE

$$J(\underline{\theta}) = \frac{1}{m} \sum_j (y^{(j)} - \underline{\theta} \cdot \underline{x}^{(j)T})^2$$

$\nabla J = ?$

$$J(\underline{\theta}) = \frac{1}{m} \sum_j \overbrace{(y^{(j)} - \theta_0 \underline{x}_0^{(j)} - \theta_1 \underline{x}_1^{(j)} - \dots)}^{e_j(\theta)}^2$$

$$\begin{aligned} \nabla J(\underline{\theta}) &= \begin{bmatrix} \frac{\partial J}{\partial \theta_0} & \frac{\partial J}{\partial \theta_1} & \dots \end{bmatrix} \\ &= \begin{bmatrix} \frac{2}{m} \sum_j -e_j(\theta) x_0^{(j)} & \frac{2}{m} \sum_j -e_j(\theta) x_1^{(j)} & \dots \end{bmatrix} \end{aligned}$$

# Gradient descent

---

Initialization

Initialize  $\theta$

Step size

Do {

- Can change as a function of iteration

$\theta \leftarrow \theta - \alpha \nabla_{\theta} J(\theta)$

} while  $(\alpha \|\nabla_{\theta} J\| > \epsilon)$

Gradient direction

Stopping condition

$$J(\underline{\theta}) = \frac{1}{m} \sum_j (y^{(j)} - \underline{\theta} \cdot \underline{x}^{(j)T})^2$$

$$\nabla J(\underline{\theta}) = -\frac{2}{m} \sum_j \underbrace{(y^{(j)} - \underline{\theta} \cdot \underline{x}^{(j)T})}_{\text{Error magnitude \& direction for datum } j} \cdot \underbrace{[x_0^{(j)} \ x_1^{(j)} \ \dots]}_{\text{Sensitivity to each } \theta_i}$$

# Derivative of MSE

Rewrite using matrix form

$$\nabla J(\underline{\theta}) = -\frac{2}{m} \sum_j \underbrace{(y^{(j)} - \underline{\theta} \cdot \underline{x}^{(j)T})}_{\text{Error magnitude \& direction for datum } j} \cdot \underbrace{[x_0^{(j)} x_1^{(j)} \dots]}_{\text{Sensitivity to each } \theta_i}$$

$$\underline{\theta} = [\theta_0, \dots, \theta_n]$$

$$\underline{y} = [y^{(1)} \dots, y^{(m)}]^T$$

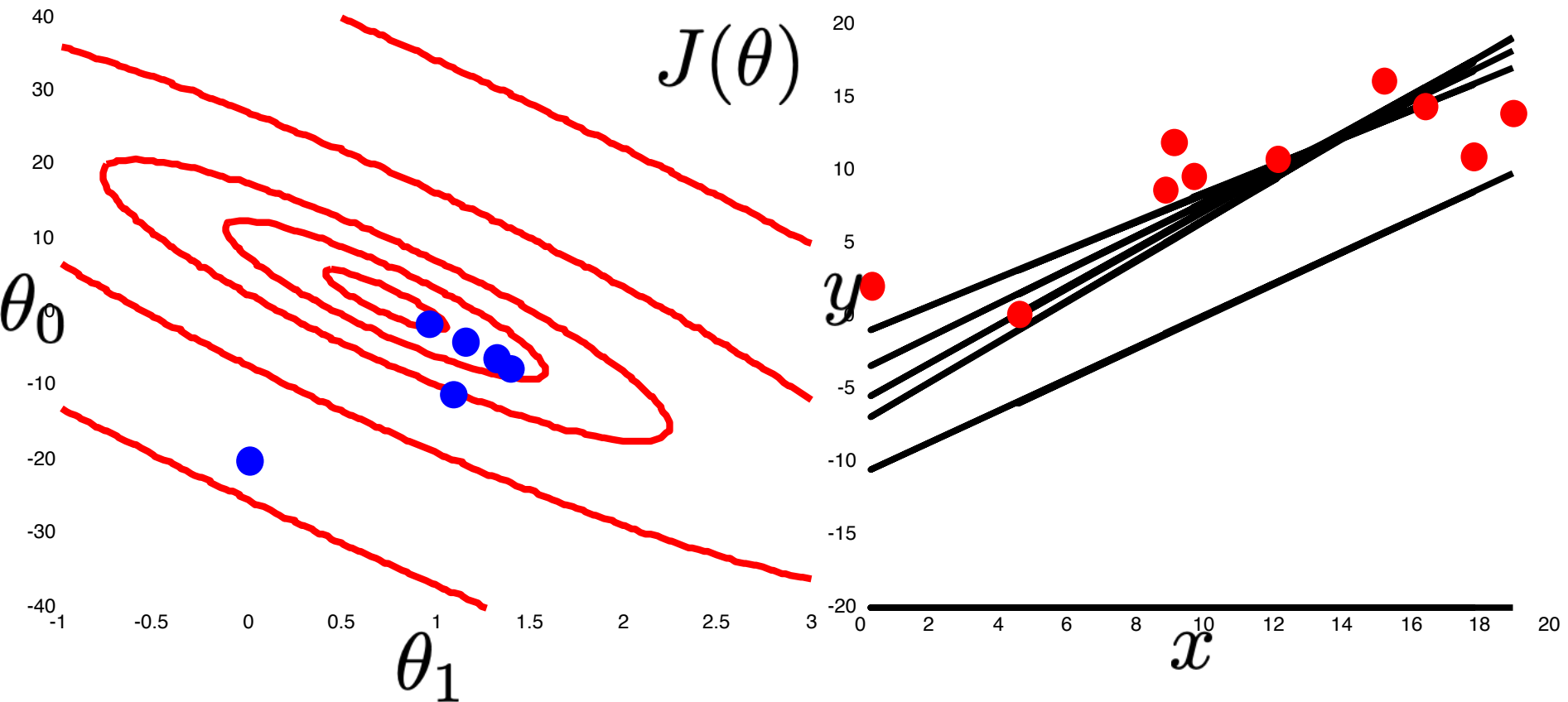
$$\nabla J(\underline{\theta}) = -\frac{2}{m} (\underline{y}^T - \underline{\theta} \underline{X}^T) \cdot \underline{X}$$

$$\underline{X} = \begin{bmatrix} x_0^{(1)} & \dots & x_n^{(1)} \\ \vdots & \ddots & \vdots \\ x_0^{(m)} & \dots & x_n^{(m)} \end{bmatrix}$$

```
e = Y - X.dot( theta.T ); # error residual
DJ = - e.dot(X) * 2.0/m   # compute the gradient
theta -= alpha * DJ       # take a step
```

# Gradient descent on cost function

---



# Comments on Gradient Descent

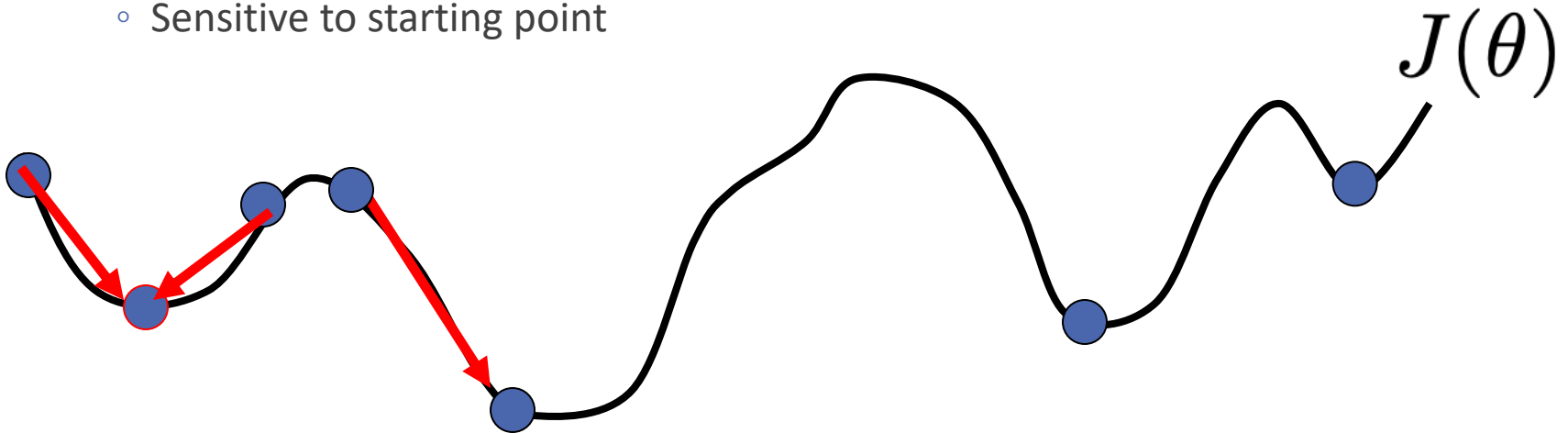
---

Very general algorithm

- We'll see it many times

Local minima

- Sensitive to starting point



# Comments on Gradient Descent

---

Very general algorithm

- We'll see it many times

Local minima

- Sensitive to starting point

Step size

- Too large? Too small? Automatic ways to choose?
- May want step size to decrease with iteration
- Common choices:
  - Fixed
  - Linear:  $C/(\text{iteration})$
  - Line search
  - Newton's method

