

# Търсене и извличане на информация. Приложение на дълбоко машинно обучение

---

Стоян Михов



Лекция 9: Намиране на градиент чрез пропагиране назад. Стохастично спускане по градиента.

# План на лекцията

---

- 1. Формалности за курса (5 мин)**
2. Намиране на градиент чрез пропагиране назад — Backpropagation (20 мин)
3. Пропагиране назад при логистична регресия (20 мин)
4. Сходимость на спускането по градиента (20 мин)
5. Стохастичен градиент (20 мин)

# Формалности

---

- Засега ще провеждаме занятията онлайн всяка сряда от 8:15 до 12:00 часа.
- Засега ще използваме платформата Google meet:  
[meet.google.com/hue-frfx-axb](https://meet.google.com/hue-frfx-axb)
- Днес ще използваме едновременно слайдове и бяла дъска. Моля следете съответния екран.
- Благодаря за предадените домашни. Ще се постараем да ги оценим до следващото занятие.
- Второто домашно задание ще бъде публикувано в Moodle около средата на декември.
- Деветата лекция се базира на глави 4 и 5 от втория учебник.

# Защо да изучаваме автоматично диференциране, спускане по градиент, стохастичен градиент и т.н.

---

- Нали в модерните системи за дълбоко обучение тези функции вече са имплементирани за нас?
- Също, защо трябва да изучаваме компилатори, след като те вече са имплементирани за нас?
  1. Да знаете какво става под повърхността винаги е полезно.
  2. Автоматичното диференциране не винаги работи перфектно — разбирането на принципите е критично при дебъгване и подобряване на моделите.
  3. При по-специални модели може да се наложи добавянето на нови модули. За разширяването на системите е съществено познаването на теорията.
  4. Може да ви се наложи да участвате в разработването на нови системи за дълбоко обучение.

# План на лекцията

---

1. Формалности за курса (5 мин)
- 2. Намиране на градиент чрез пропагиране назад — Backpropagation (20 мин)**
3. Пропагиране назад при логистична регресия (20 мин)
4. Сходимость на спускането по градиента (20 мин)
5. Стохастичен градиент (20 мин)

# Числено намиране на градиент

---

- Нека  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ . Тогава:  $\frac{\partial f(x_1, \dots, x_n)}{\partial x_i} = \lim_{h \rightarrow 0} \frac{f(x_1, \dots, x_i + h, \dots, x_n) - f(x_1, \dots, x_i, \dots, x_n)}{h}$  т.е.  
 $\frac{\partial f(\mathbf{x})}{\partial x_i} = \lim_{h \rightarrow 0} \frac{f(\mathbf{x} + h\mathbf{e}_i) - f(\mathbf{x})}{h}$ , където  $\mathbf{e}_i$  е  $i$ -тия базисен вектор.
- Полагайки достатъчно малко  $h$ , примерно  $h = 10^{-4}$ , ние можем да намерим приближение на производната:  $\frac{\partial f(\mathbf{x})}{\partial x_i} \approx \frac{f(\mathbf{x} + h\mathbf{e}_i) - f(\mathbf{x})}{h}$ .
- По добра апроксимация на производната:  $\frac{\partial f(\mathbf{x})}{\partial x_i} \approx \frac{f(\mathbf{x} + h\mathbf{e}_i) - f(\mathbf{x} - h\mathbf{e}_i)}{2h}$ .
- Градиента намираме, като апроксимираме производната в дадената точка по всяка от  $n$ -те координати.
- **Задача:** Докажете, че втората формула ни дава по-добро приближение на производната, като оцените порядъка на грешката.
- **Сложност:** Ако сложността на израза за  $f$  е от порядък  $O(k)$  то сложността за численото намиране на градиента на  $f$  в точката  $\mathbf{x}$  е от порядък  $O(nk)$ .

# Аналитично намиране на градиент

---

- Нека  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ . Нека сме намерили аналитични изрази за производните  $\frac{\partial f(\mathbf{x})}{\partial x_i}$  за  $i = 1, 2, \dots, n$ .
- Производната по дадено направление намираме, като заместим в съответния израз за производната с дадената точка.
- Градиента намираме, като заместим в израза за производната по всяка от  $n$ -те координати с дадената точка.
- **Сложност:** Ако сложността на изразите за  $\frac{\partial f(\mathbf{x})}{\partial x_i}$  е от порядък  $O(k')$ , то сложността за аналитичното намиране на градиента на  $f$  в точката  $\mathbf{x}$  е от порядък  $O(nk')$ .
- **Задача:** Докажете, че съществуват изрази със сложност  $O(k)$ , за които сложността на израза за производната е от порядък  $O(2^k)$ .

# Намиране на градиент чрез пропагиране назад

## — **Backpropagation**

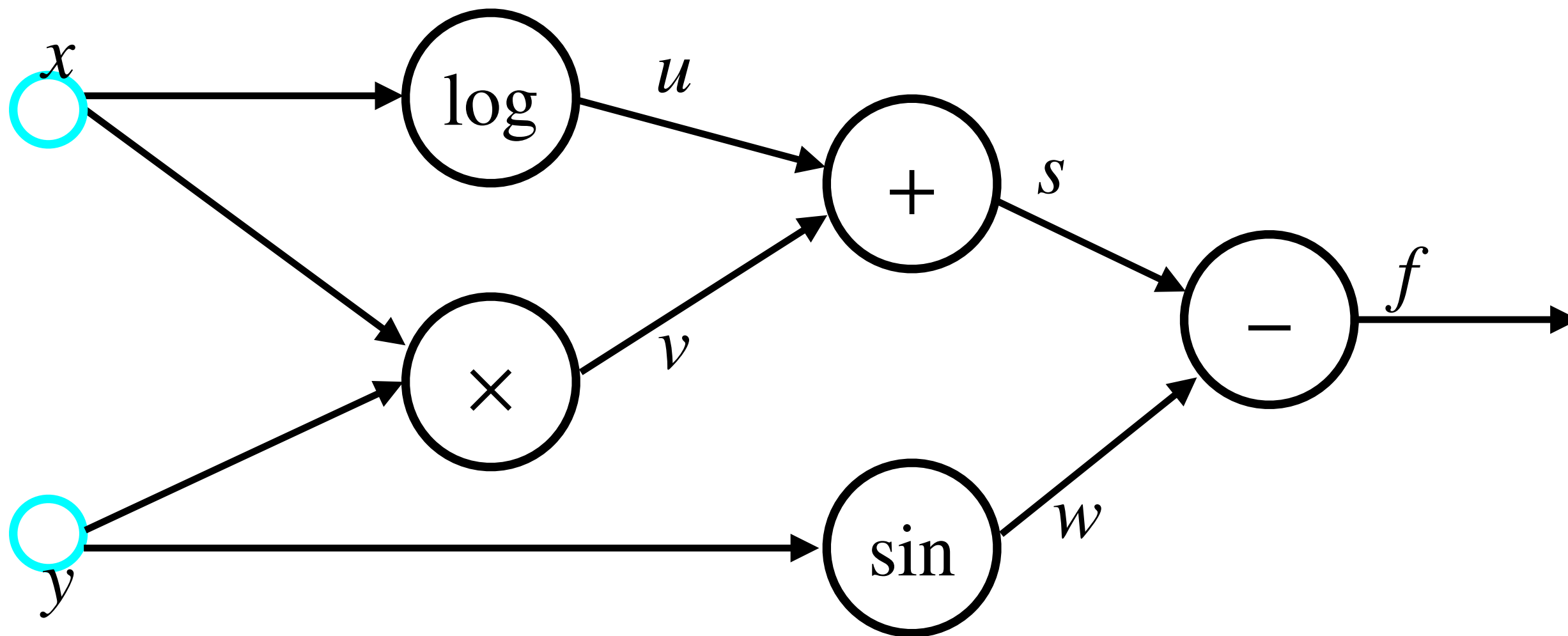
---

- Намира точните производни и градиент в дадена точка — не е числено приближение.
- Използва се аналитичен израз за производните само за локалните функции — избягва се намирането на изразите за производните на целевата функция.
- Преизползват се междинните резултати от изчисленията, с което се постига оптимална изчислителна сложност.
- Сложността за намиране на градиента е  $O(k)$ , за израз със сложност  $k$ .



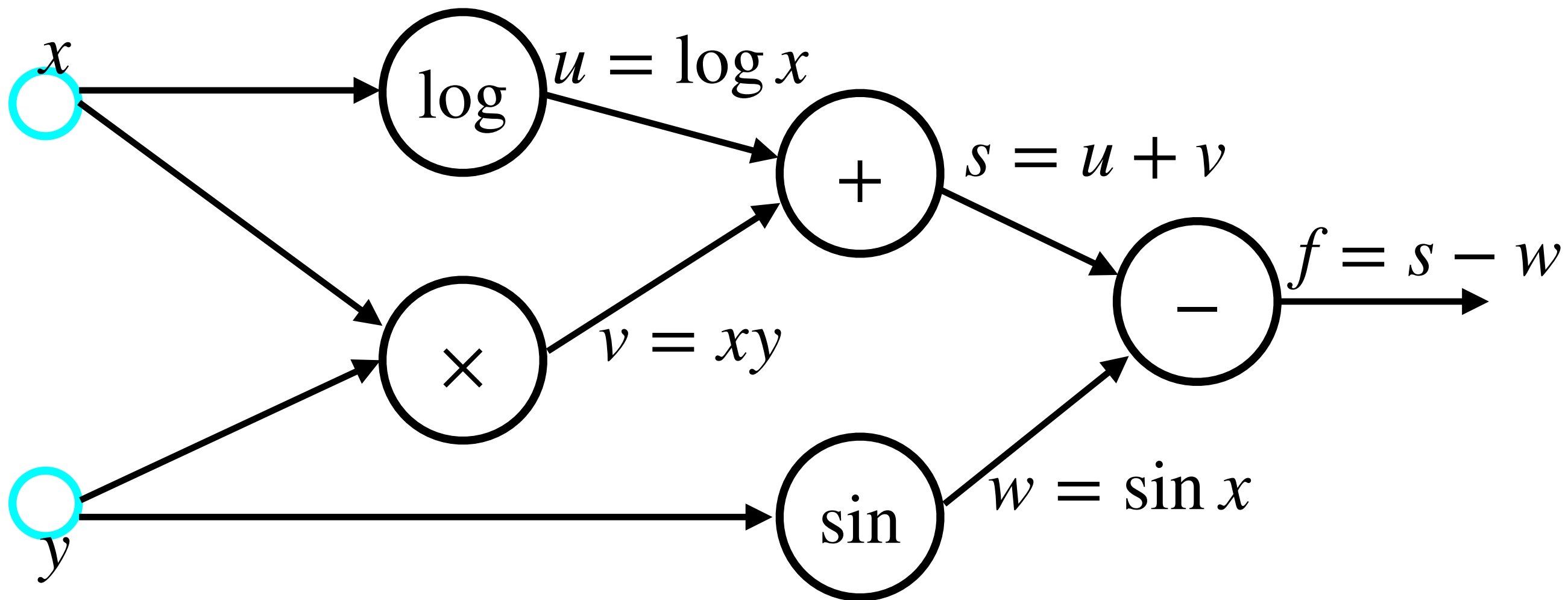
Пример:  $f(x, y) = (\log(x) + xy) - \sin(y)$

---



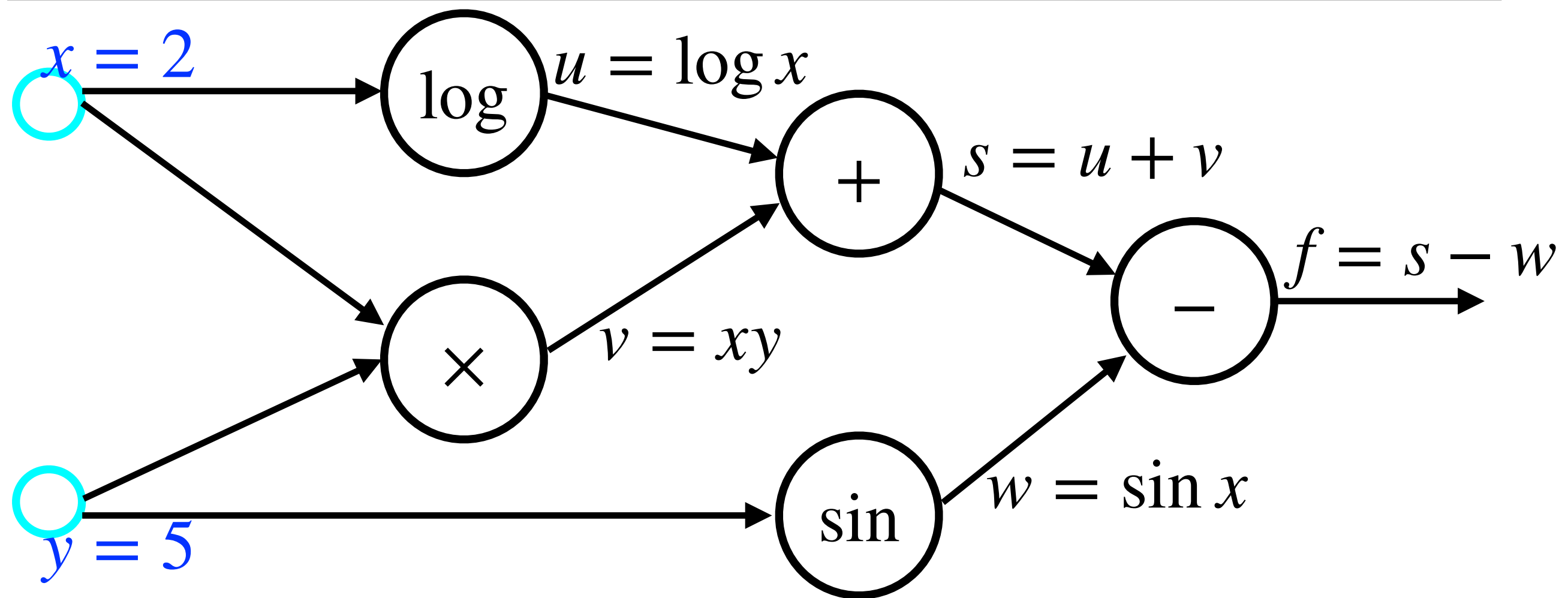
Пример:  $f(x, y) = (\log(x) + xy) - \sin(y)$

---



Пример:  $f(x, y) = (\log(x) + xy) - \sin(y)$

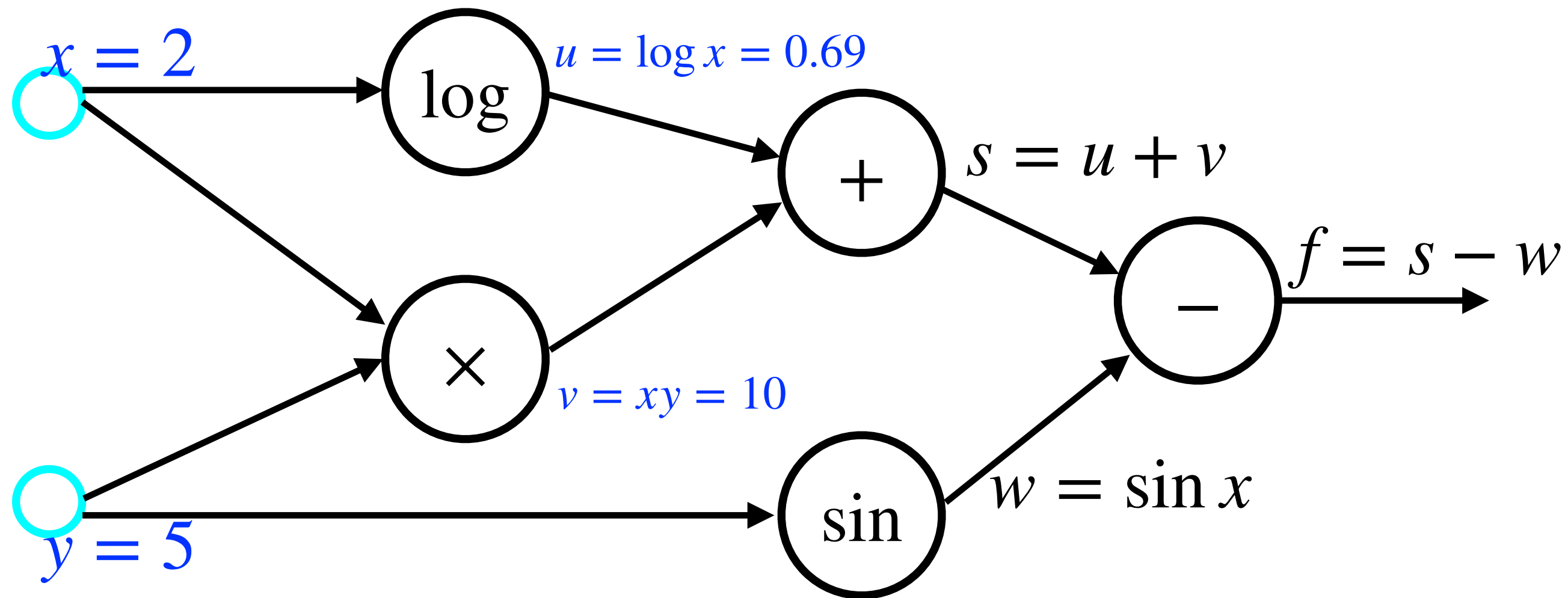
---



**Пропагиране напред — Forward propagation**

Пример:  $f(x, y) = (\log(x) + xy) - \sin(y)$

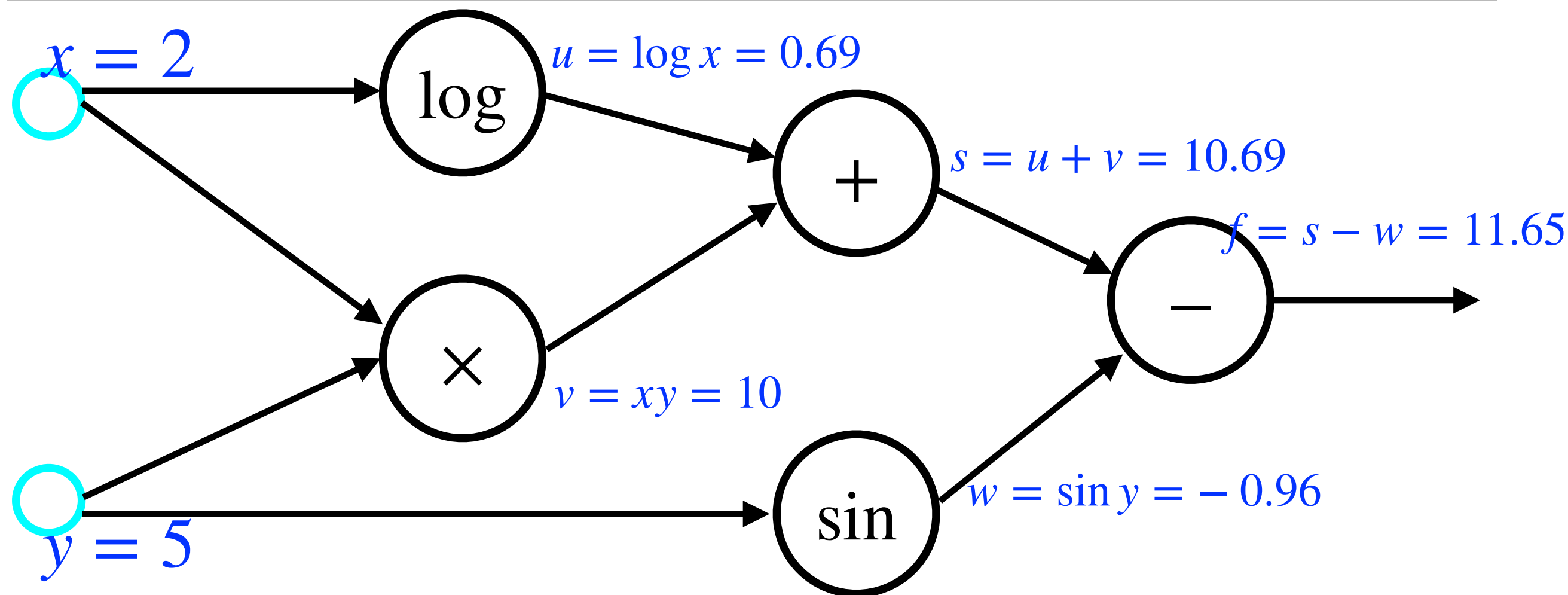
---



**Пропагиране напред — Forward propagation**

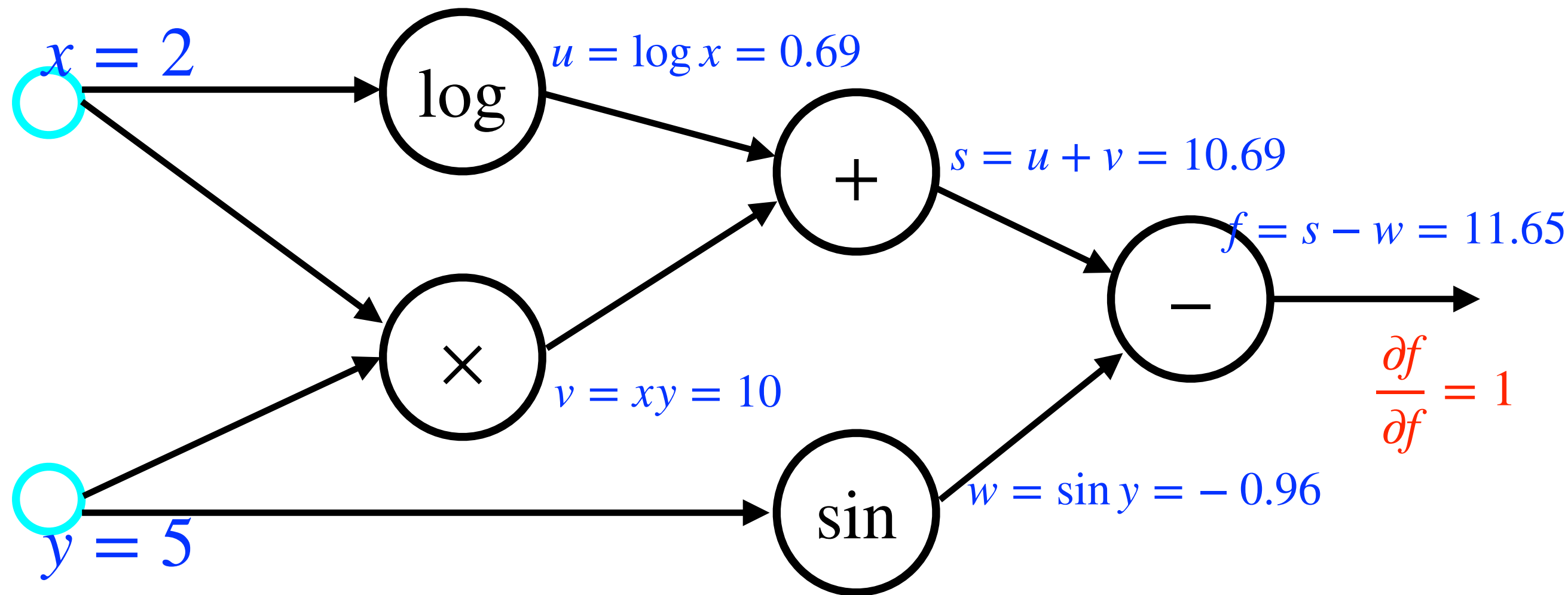
Пример:  $f(x, y) = (\log(x) + xy) - \sin(y)$

---



**Пропагиране напред — Forward propagation**

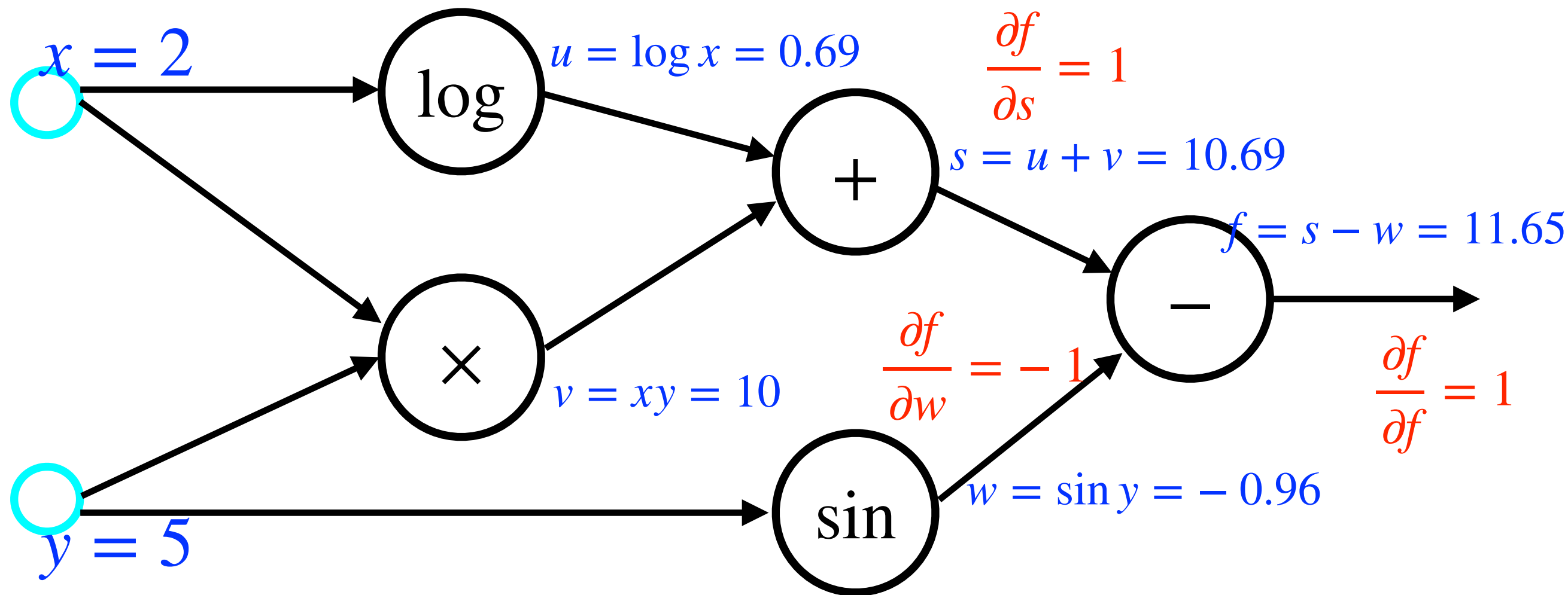
Пример:  $f(x, y) = (\log(x) + xy) - \sin(y)$



**Пропагиране напред — Forward propagation**

**Пропагиране назад — Back propagation**

Пример:  $f(x, y) = (\log(x) + xy) - \sin(y)$

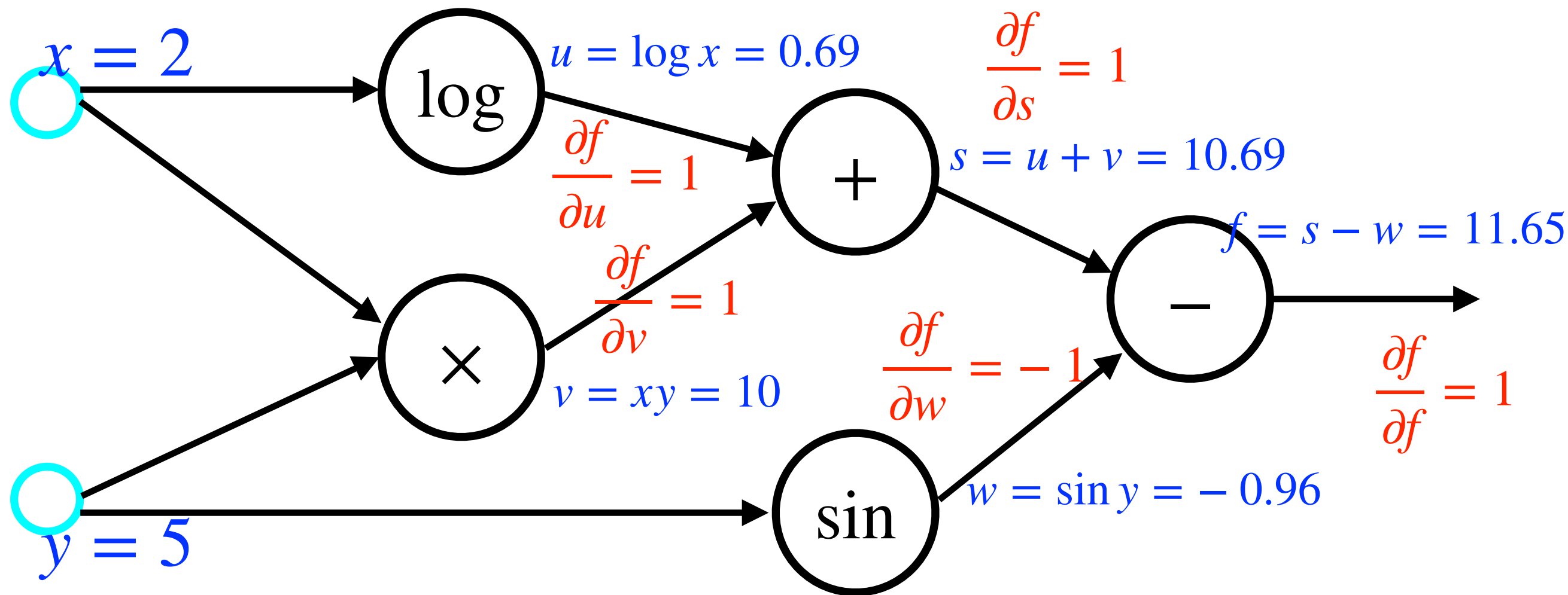


**Пропагиране напред — Forward propagation**

**Пропагиране назад — Back propagation**

$$\frac{\partial f}{\partial w} = \frac{\partial f}{\partial f} \frac{\partial f}{\partial w} = \frac{\partial f}{\partial f} (-1) = -1, \quad \frac{\partial f}{\partial s} = \frac{\partial f}{\partial f} \frac{\partial f}{\partial s} = \frac{\partial f}{\partial f} 1 = 1$$

Пример:  $f(x, y) = (\log(x) + xy) - \sin(y)$



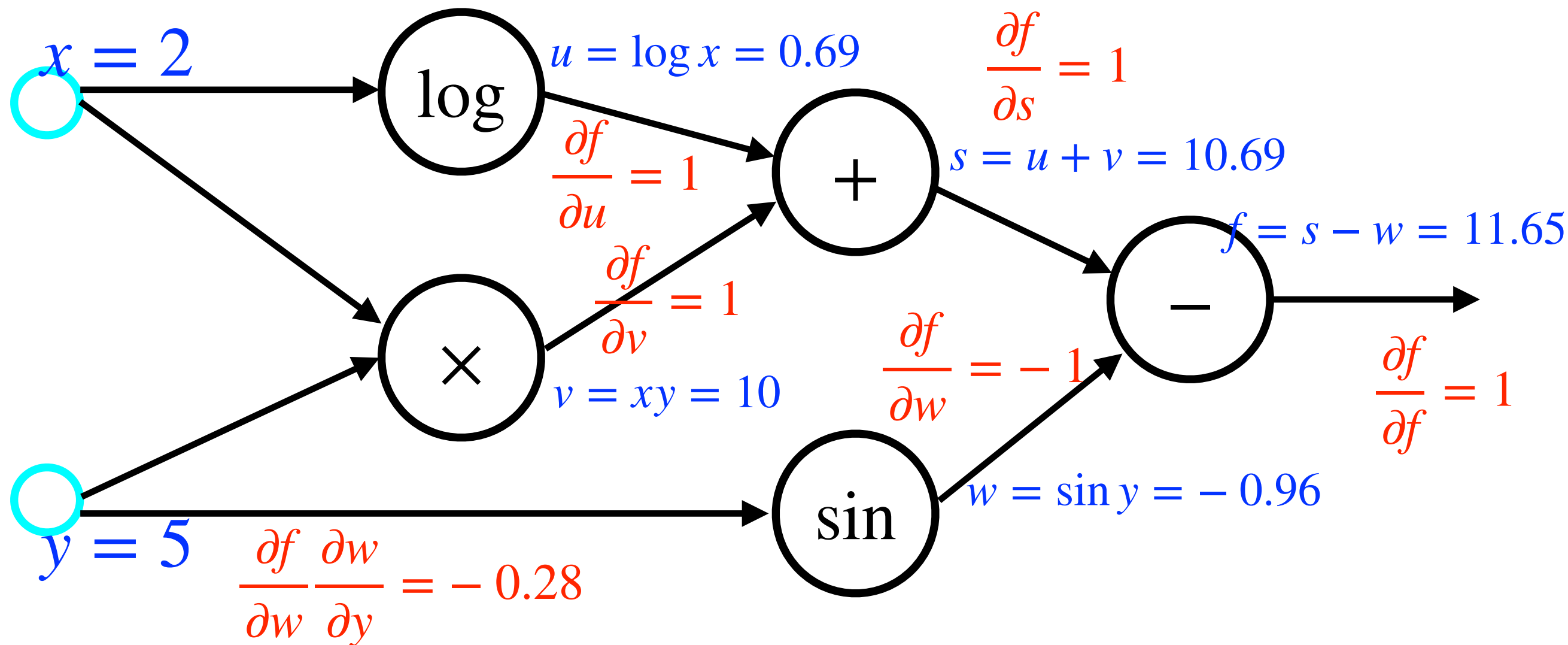
**Пропагиране напред — Forward propagation**

**Пропагиране назад — Back propagation**

$$\frac{\partial f}{\partial u} = \frac{\partial f}{\partial s} \frac{\partial s}{\partial u} = \frac{\partial f}{\partial s} 1 = 1, \quad \frac{\partial f}{\partial v} = \frac{\partial f}{\partial s} \frac{\partial s}{\partial v} = \frac{\partial f}{\partial s} 1 = 1$$



Пример:  $f(x, y) = (\log(x) + xy) - \sin(y)$

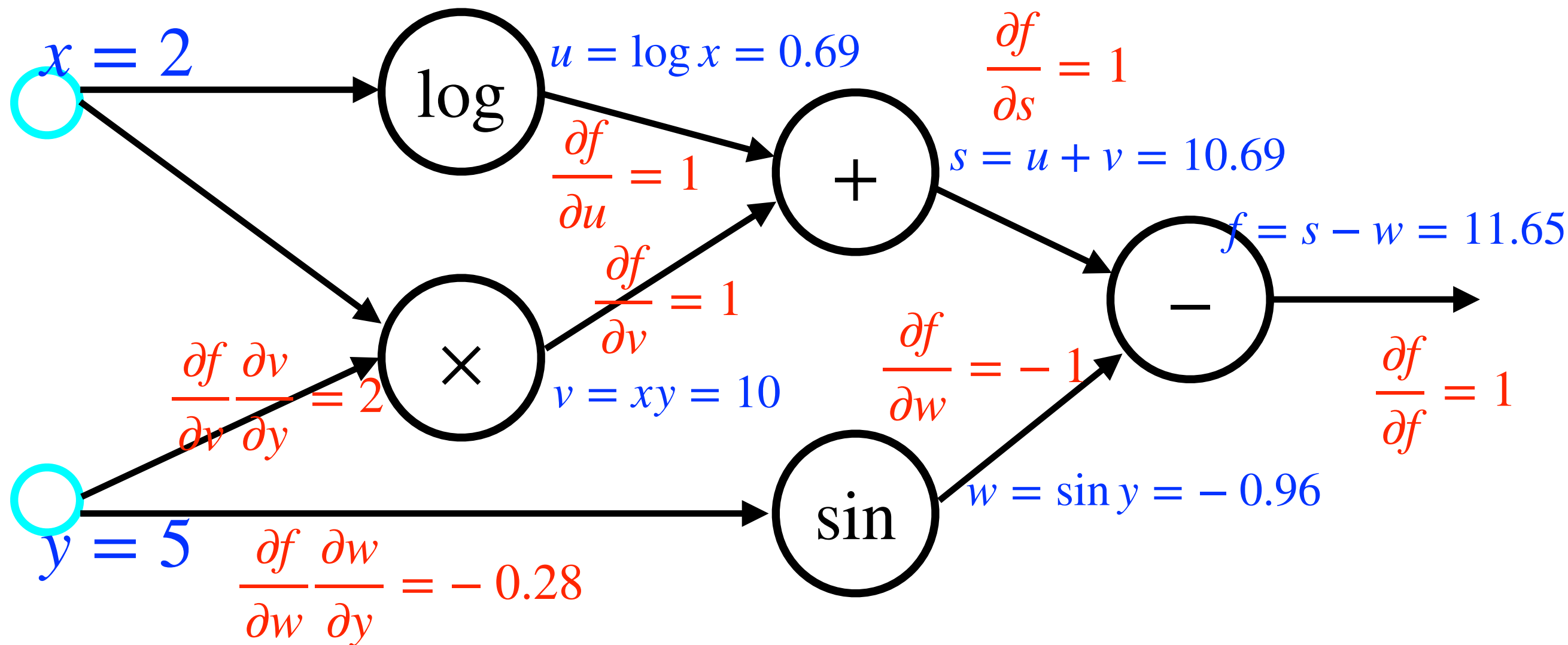


**Пропагиране напред — Forward propagation**

**Пропагиране назад — Back propagation**

$$\frac{\partial f}{\partial w} \frac{\partial w}{\partial y} = \frac{\partial f}{\partial w} \cos(y) = (-1) \cos(5) = -0.28$$

Пример:  $f(x, y) = (\log(x) + xy) - \sin(y)$

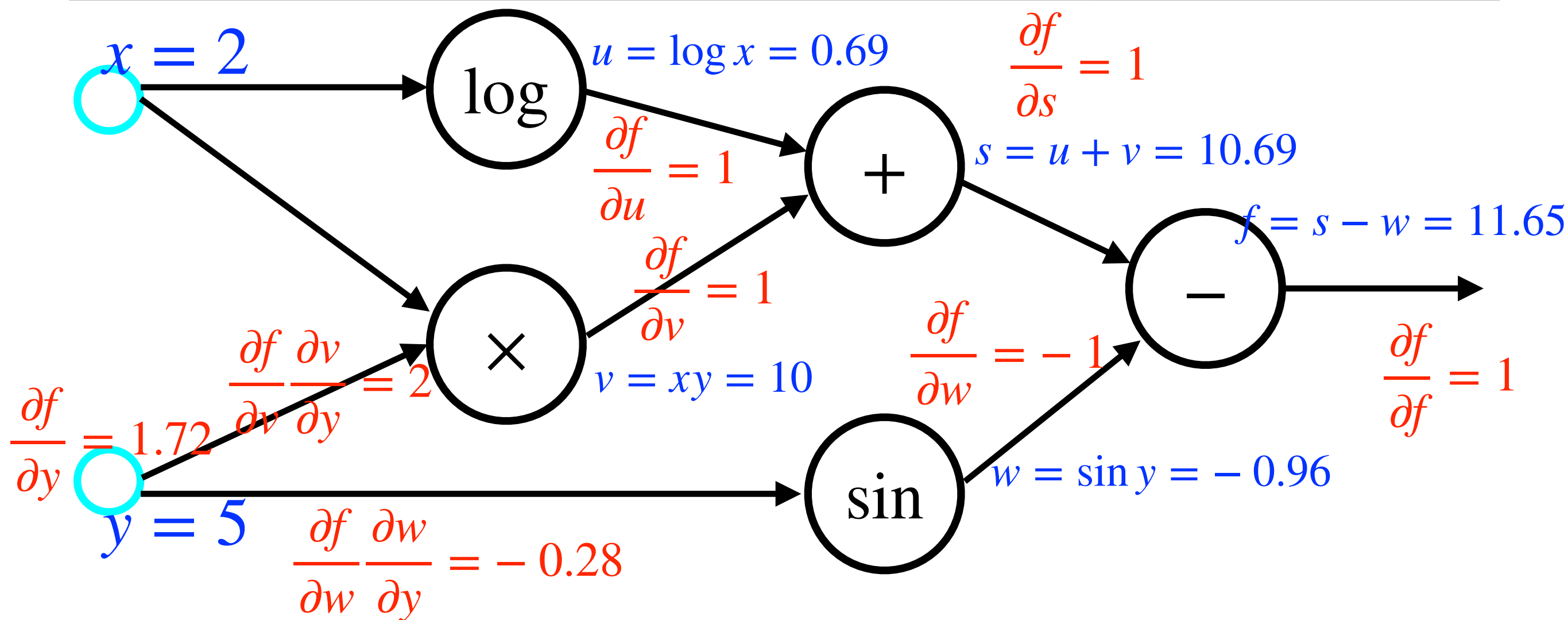


**Пробагиране напред — Forward propagation**

**Пробагиране назад — Back propagation**

$$\frac{\partial f}{\partial v} \frac{\partial v}{\partial y} = \frac{\partial f}{\partial v} x = 2$$

Пример:  $f(x, y) = (\log(x) + xy) - \sin(y)$

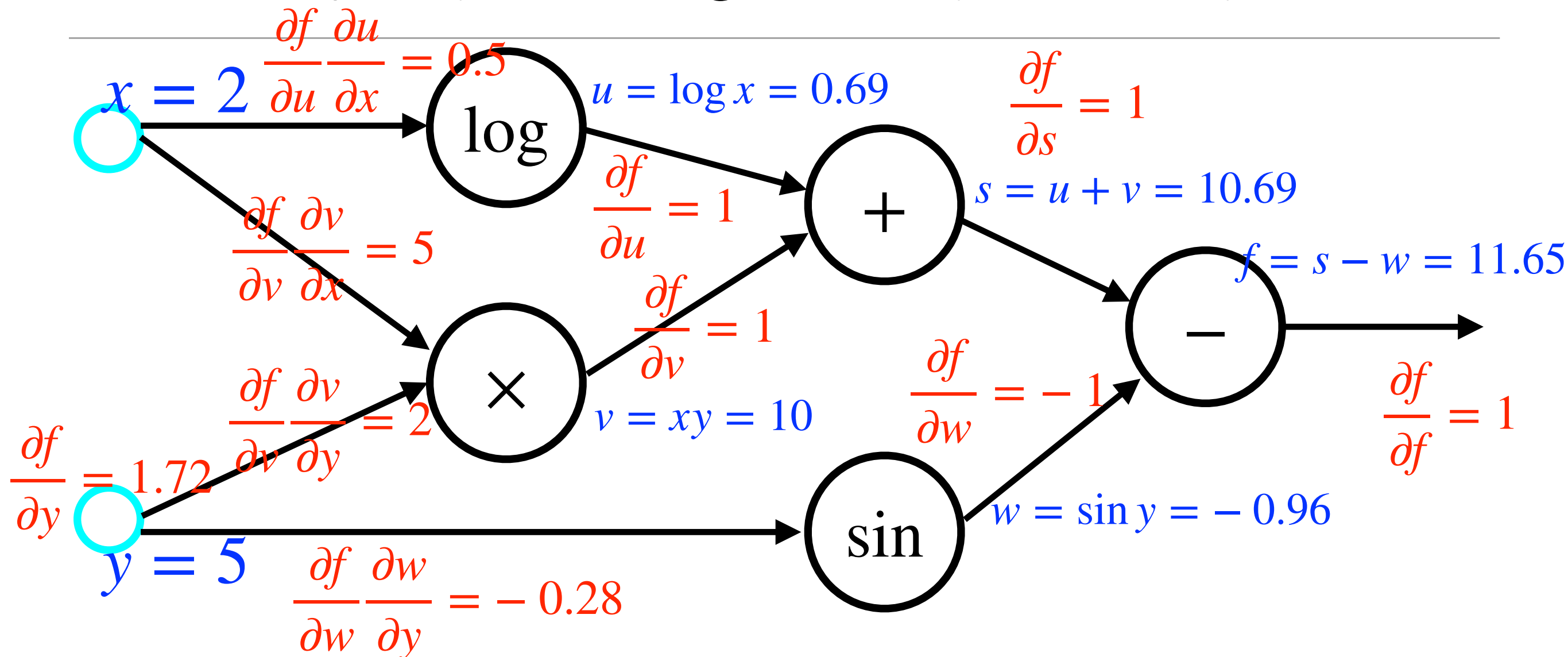


**Пропагиране напред — Forward propagation**

**Пропагиране назад — Back propagation**

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial v} \frac{\partial v}{\partial y} + \frac{\partial f}{\partial w} \frac{\partial w}{\partial y} = 2 - 0.28 = 1.72$$

Пример:  $f(x, y) = (\log(x) + xy) - \sin(y)$

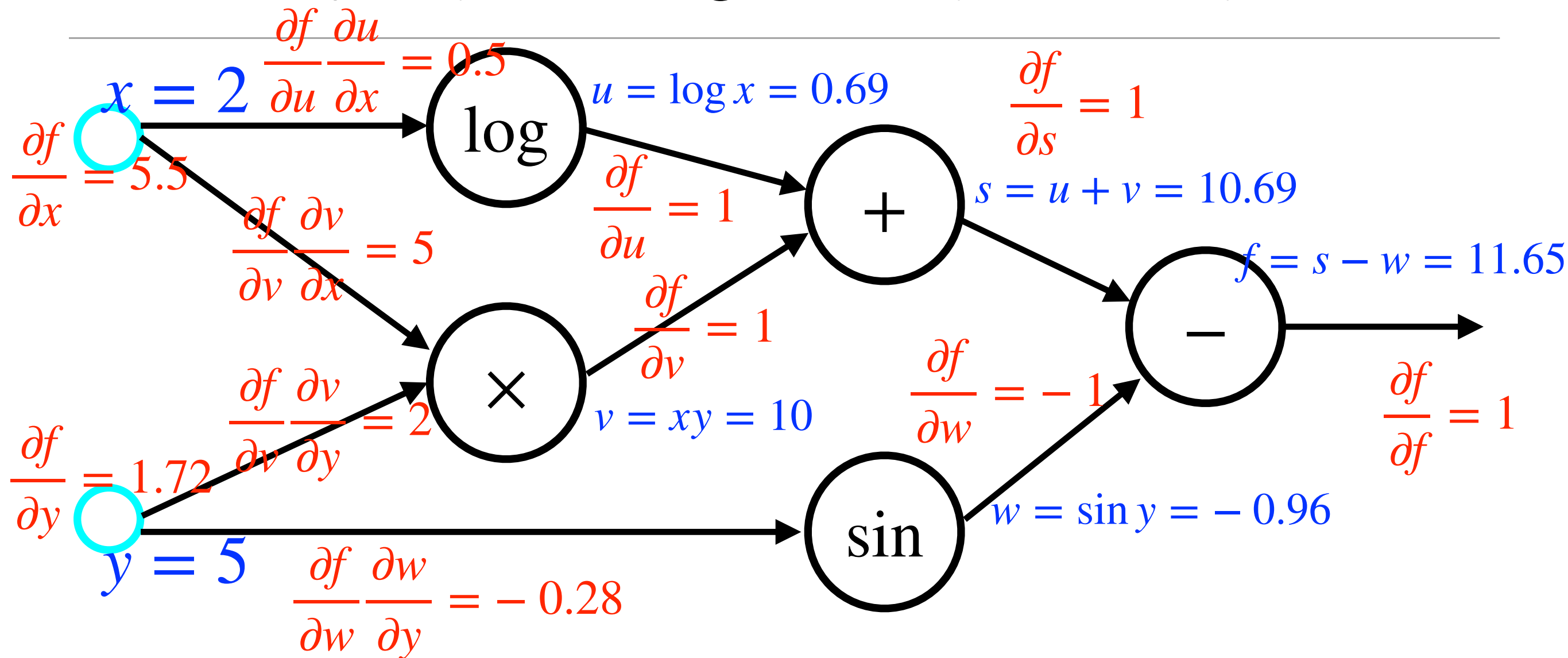


**Пропагиране напред — Forward propagation**

**Пропагиране назад — Back propagation**

$$\frac{\partial f}{\partial u} \frac{\partial u}{\partial x} = \frac{\partial f}{\partial u} \frac{1}{x} = 0.5, \quad \frac{\partial f}{\partial v} \frac{\partial v}{\partial x} = \frac{\partial f}{\partial v} y = 5$$

Пример:  $f(x, y) = (\log(x) + xy) - \sin(y)$



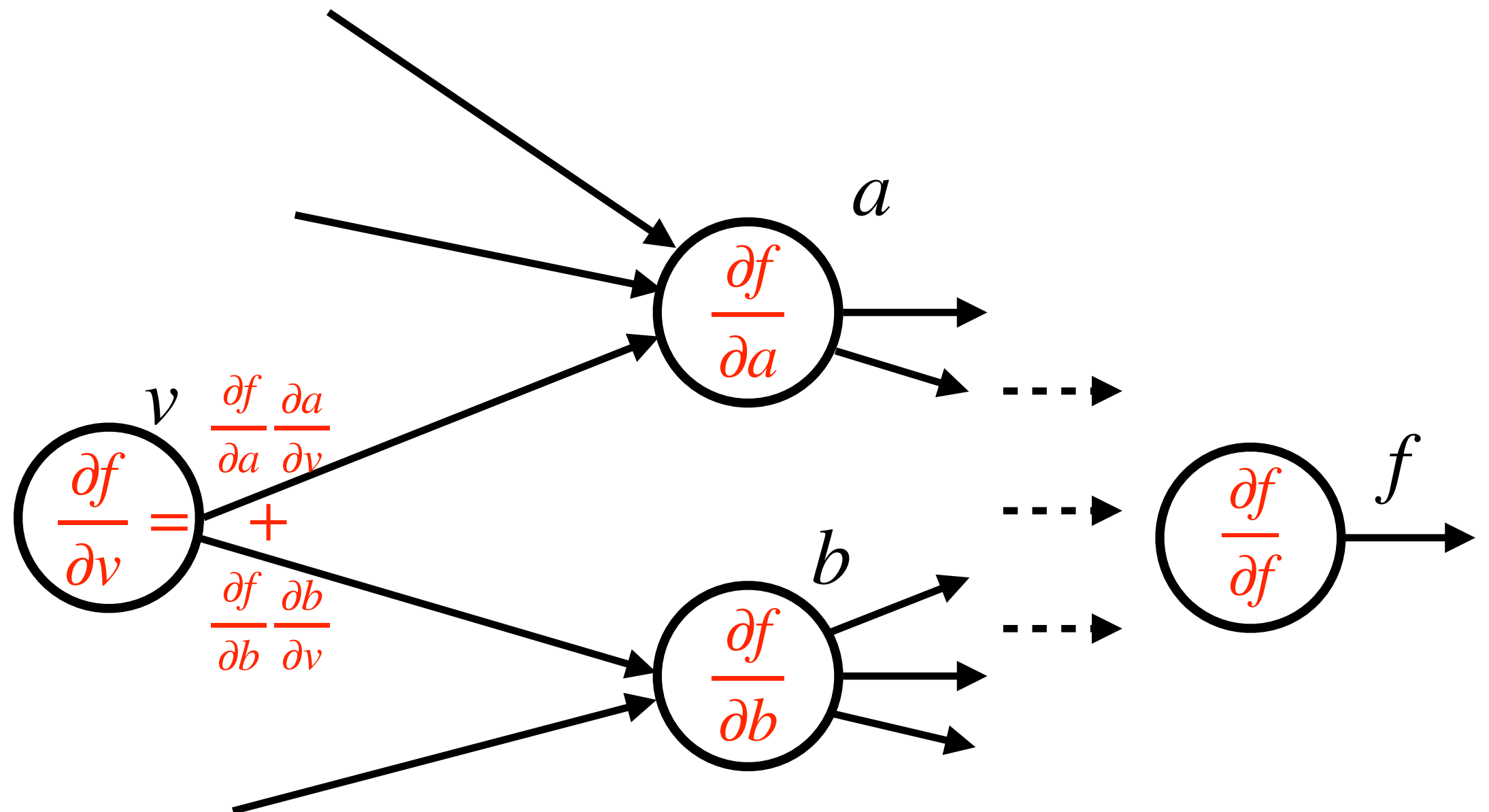
**Пропагиране напред — Forward propagation**

**Пропагиране назад — Back propagation**

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial u} \frac{\partial u}{\partial x} + \frac{\partial f}{\partial v} \frac{\partial v}{\partial x} = 5 + 0.5 = 5.5$$

# Основен инвариант

---



# Формализация на Backpropagation — изчислителен граф

---

- Даден е ацикличен граф  $G = (V, E), E \subset V \times V$ .
- За всеки връх  $v \in V$  (непосредствените) предшественици означаваме с  $P_G(v) = [p \in V \mid (p, v) \in E]$ . Ще предпологаеме, че  $P_G(v)$  е списък.
- Листата на графа ще означаваме с  $L(G) = [v \mid P_G(v) = \emptyset]$ .
- Крайни върхове на графа ще означаваме с  $T(G) = [v \mid \neg \exists u \in V : (v, u) \in E]$ . Ще предпологаеме, че в графа съществува единствен краен връх:  $|T(G)| = 1$

# Формализация на Backpropagation — изчислителен граф

---

- За всеки връх  $v \in V \setminus L(G)$  ще предпологаме, че е дефинирана функция за изчисляване на стойността:  
 $\text{calc}(G, v) : \mathbb{R}^{|P_G(v)|} \rightarrow \mathbb{R}$ .
- За всеки връх  $v \in V \setminus L(G)$  и всеки негов предшественик  $u \in P_G(v)$  ще предпологаме, че е дефинирана функция за изчисляване на частната производна на функцията  $\text{calc}(G, v)$  спрямо аргумента  $u$ , която бележим с  
 $\text{deriv}(G, v, u) : \mathbb{R}^{|P_G(v)|} \rightarrow \mathbb{R}$ .



# Backpropagation алгоритъм

---

Forward(G):

```
1  for v in topologicalSort(G) do
2      if v in Leafs(G) then read(F(v))
3      else
4          [p1,p2,...,pk] <- Predecessors(G)(v)
5          F(v) <- calc(G,v)(F(p1),...,F(pk))
6  return F
```

Backward(G,F):

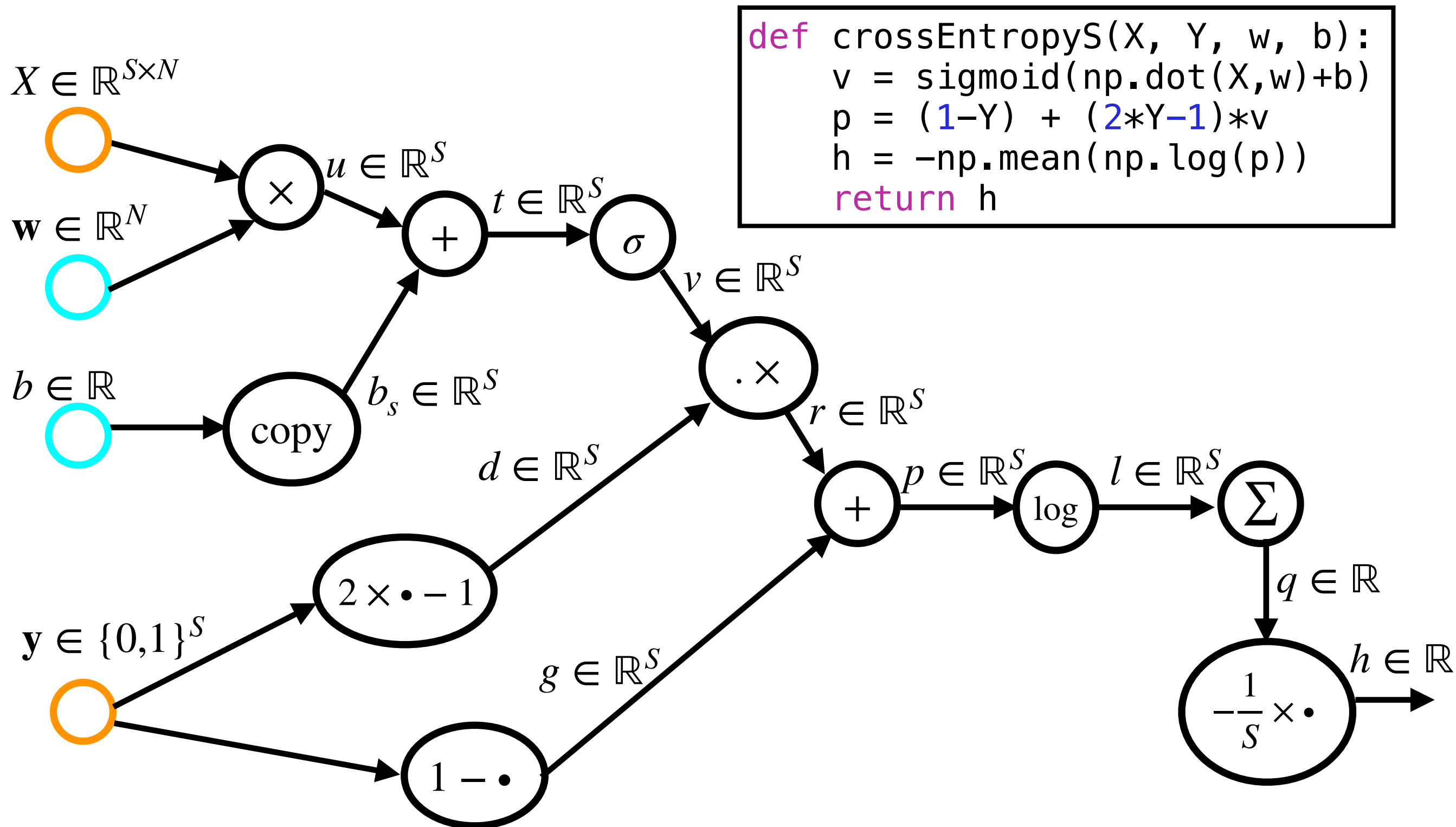
```
1  for v in topologicalSort(G) do B(v) <- 0
2  for v in reverseTopologicalSort(G) do
3      if v in Top(G) then B(v) <- 1
4      [p1,p2,...,pk] <- Predecessors(G)(v)
5      for p in Predecessors(G)(v) do
6          B(p) += B(v) * deriv(G,v,p)(F(p1),...,F(pk))
7  return B
```

# План на лекцията

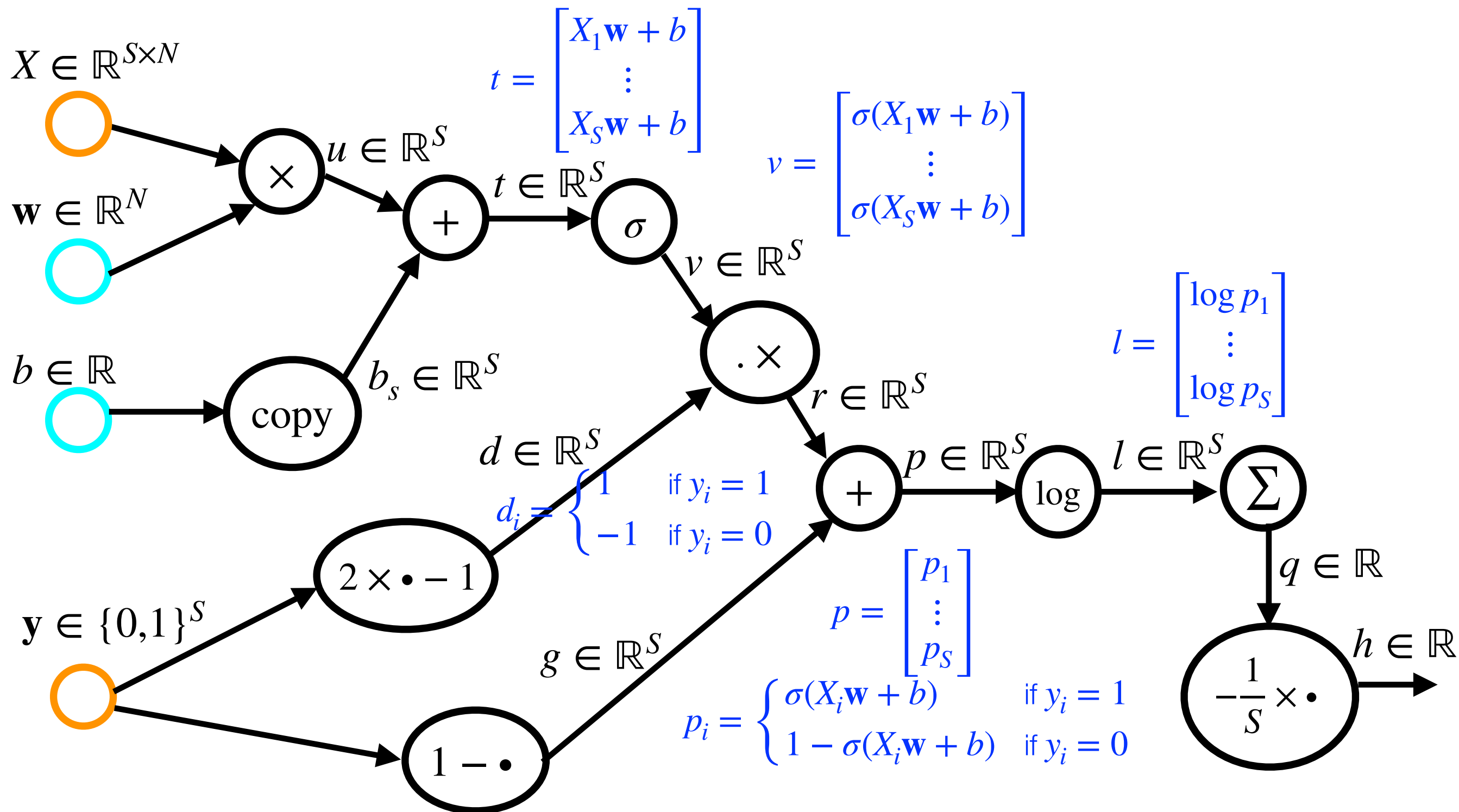
---

1. Формалности за курса (5 мин)
2. Намиране на градиент чрез пропагиране назад — Backpropagation (20 мин)
- 3. Пропагиране назад при логистична регресия (20 мин)**
4. Сходимость на спускането по градиента (20 мин)
5. Стохастичен градиент (20 мин)

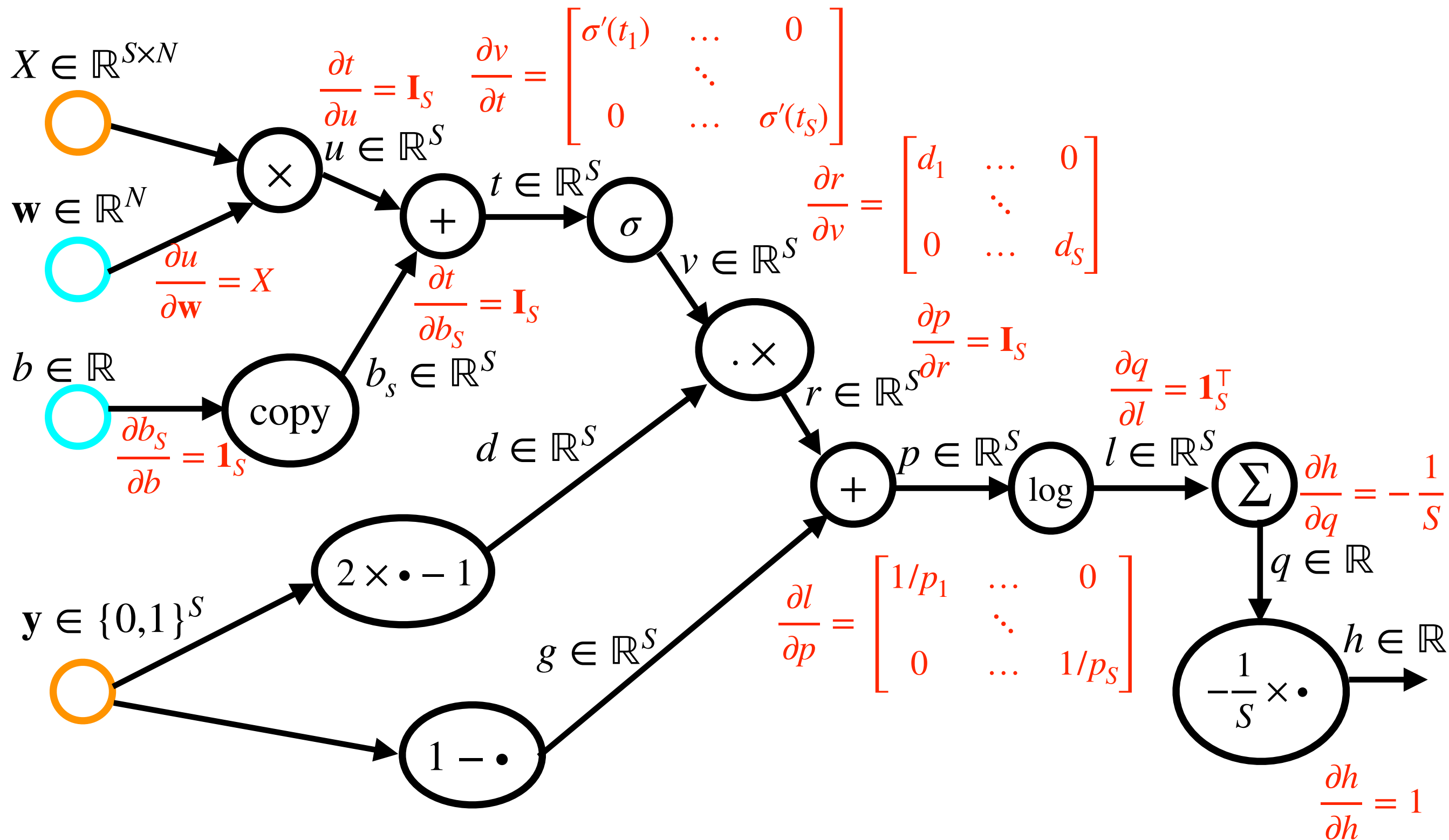
# Логистична регресия — векторен запис



# Логистична регресия — векторен запис



# Логистична регресия — векторен запис



$$\begin{aligned}
\frac{\partial h}{\partial t} &= \frac{\partial h}{\partial q} \frac{\partial q}{\partial l} \frac{\partial l}{\partial p} \frac{\partial p}{\partial r} \frac{\partial r}{\partial v} \frac{\partial v}{\partial t} \\
&= -\frac{1}{S} \mathbf{1}_S^\top \begin{bmatrix} 1/p_1 & \dots & 0 \\ & \ddots & \\ 0 & \dots & 1/p_S \end{bmatrix} \mathbf{I}_S \begin{bmatrix} d_1 & \dots & 0 \\ & \ddots & \\ 0 & \dots & d_S \end{bmatrix} \begin{bmatrix} \sigma'(t_1) & \dots & 0 \\ & \ddots & \\ 0 & \dots & \sigma'(t_S) \end{bmatrix} = \\
&= \left[ -\frac{y_1 - \sigma(X_1 \mathbf{w} + b)}{S} \dots -\frac{y_S - \sigma(X_S \mathbf{w} + b)}{S} \right]
\end{aligned}$$

$$\frac{\partial h}{\partial b} = \frac{\partial h}{\partial t} \frac{\partial t}{\partial b_S} \frac{\partial b_S}{\partial b} = \left[ -\frac{y_1 - \sigma(X_1 \mathbf{w} + b)}{S} \quad \dots \quad -\frac{y_S - \sigma(X_S \mathbf{w} + b)}{S} \right] \mathbf{I}_S \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} = -\frac{1}{S} \sum_{i=1}^S y_i - \sigma(X_i \mathbf{w} + b)$$

$$\begin{aligned}
\frac{\partial h}{\partial \mathbf{w}} &= \frac{\partial h}{\partial t} \frac{\partial t}{\partial u} \frac{\partial u}{\partial \mathbf{w}} = \left[ -\frac{y_1 - \sigma(X_1 \mathbf{w} + b)}{S} \dots -\frac{y_S - \sigma(X_S \mathbf{w} + b)}{S} \right] \mathbf{I}_S X = \\
&= -\frac{1}{S} \sum_{i=1}^S (y_i - \sigma(X_i \mathbf{w} + b)) X_i
\end{aligned}$$

# План на лекцията

---

1. Формалности за курса (5 мин)
2. Намиране на градиент чрез пропагиране назад — Backpropagation (20 мин)
3. Пропагиране назад при логистична регресия (20 мин)
4. **Сходимость на спускането по градиента (20 мин)**
5. Стохастичен градиент (20 мин)

# Градиент, Хесиан, развиване в ред на Тейлър

---

**Теорема (Тейлър):** Нека  $g : \mathbb{R} \rightarrow \mathbb{R}$  е два пъти диференцируема с непрекъснати производни в околност на точката  $t_0 \in \mathbb{R}$ . Нека  $t \in \mathbb{R}$  е произволна точка в тази околност. Тогава съществува  $\bar{t} \in (t_0, t)$ , така че:

$$g(t) = g(t_0) + g'(t_0)(t - t_0) + \frac{1}{2}g''(\bar{t})(t - t_0)^2$$

Ще изведем теоремата на Тейлър за многомерния случай.

Нека  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ . **Градиент** и **Хесиан** на  $f$  наричаме съответно:

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \frac{\partial}{\partial \mathbf{x}} f(\mathbf{x}) = \begin{bmatrix} \frac{\partial}{\partial x_1} f(\mathbf{x}) \\ \frac{\partial}{\partial x_2} f(\mathbf{x}) \\ \vdots \\ \frac{\partial}{\partial x_n} f(\mathbf{x}) \end{bmatrix}, \quad \nabla_{\mathbf{x}}^2 f(\mathbf{x}) = \frac{\partial^2}{\partial \mathbf{x}^2} f(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2}{\partial x_1 \partial x_1} f(\mathbf{x}) & \cdots & \frac{\partial^2}{\partial x_1 \partial x_n} f(\mathbf{x}) \\ \vdots & \ddots & \vdots \\ \frac{\partial^2}{\partial x_n \partial x_1} f(\mathbf{x}) & \cdots & \frac{\partial^2}{\partial x_n \partial x_n} f(\mathbf{x}) \end{bmatrix}$$



Нека  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ . Т.е. за  $\mathbf{v} \in \mathbb{R}^n$  имаме, че  $f(\mathbf{v}) \in \mathbb{R}$ . Разглеждаме функциите  $\mathbf{v} : \mathbb{R} \rightarrow \mathbb{R}^n$ ,  $\mathbf{v}(t) = \mathbf{x}_0 + t(\mathbf{x} - \mathbf{x}_0)$  и  $g : \mathbb{R} \rightarrow \mathbb{R}$ ,  $g(t) = f(\mathbf{v}(t)) = f(\mathbf{x}_0 + t(\mathbf{x} - \mathbf{x}_0))$  и прилагаме теоремата на Тейлър за  $g$  при  $t_0 = 0$ ,  $t = 1$ . В такъв случай имаме, че

$$g'(t) = (f(\mathbf{v}(t)))' = \left( \frac{\partial f(\mathbf{v})}{\partial \mathbf{v}} \right)^{\top} \frac{\partial \mathbf{v}(t)}{\partial t} = (\nabla_{\mathbf{v}} f(\mathbf{v}))^{\top} (\mathbf{x} - \mathbf{x}_0),$$

$$g''(t) = \frac{\partial}{\partial t} \left( \left( \frac{\partial f(\mathbf{v})}{\partial \mathbf{v}} \right)^{\top} (\mathbf{x} - \mathbf{x}_0) \right) = (\mathbf{x} - \mathbf{x}_0)^{\top} \frac{\partial^2 f(\mathbf{v})}{\partial \mathbf{v}^2} (\mathbf{x} - \mathbf{x}_0) = (\mathbf{x} - \mathbf{x}_0)^{\top} \nabla_{\mathbf{v}}^2 f(\mathbf{v}) (\mathbf{x} - \mathbf{x}_0)$$

**Теорема (Тейлър):** Нека  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  е два пъти диференцируема с непрекъснати производни в околност на точката  $\mathbf{x}_0 \in \mathbb{R}^n$ . Нека  $\mathbf{x} \in \mathbb{R}^n$  е произволна точка в тази околност. Тогава съществува  $\bar{t} \in (0,1)$ , така че ако  $\bar{\mathbf{x}} = \mathbf{x}_0 + \bar{t}(\mathbf{x} - \mathbf{x}_0)$ , то:

$$f(\mathbf{x}) = f(\mathbf{x}_0) + (\nabla_{\mathbf{x}} f(\mathbf{x}_0))^{\top} (\mathbf{x} - \mathbf{x}_0) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_0)^{\top} \nabla_{\mathbf{x}}^2 f(\bar{\mathbf{x}}) (\mathbf{x} - \mathbf{x}_0)$$

# Сходимость на спускането по градиента

---

- Стремим се да минимизираме кросентропията

$$H_X[\text{Pr} \parallel \text{Pr}_\theta] = -\frac{1}{|X|} \sum_{i=1}^{|X|} \log \text{Pr}_\theta[\mathbf{x}^{(i)}]. \text{ Нека положим}$$

$$f(\theta) = H_X[\text{Pr} \parallel \text{Pr}_\theta].$$

- Спускаме се по градиента:  $\theta_{k+1} = \theta_k - \alpha \nabla_\theta f(\theta_k)$

- От развиването в ред на Тейлър получаваме:

$$f(\theta_{k+1}) = f(\theta_k) - \alpha \nabla f(\theta_k)^\top \nabla f(\theta_k) + \frac{1}{2} (\alpha \nabla f(\theta_k))^\top \nabla^2 f(\bar{\theta}) (\alpha \nabla f(\theta_k))$$

- Да допуснем, че вторите производни са ограничени:  
 $\|\nabla^2 f(\theta)\| \leq L$ . Тогава:  $\|\mathbf{u}^\top \nabla^2 f(\theta) \mathbf{u}\| \leq L\|\mathbf{u}\|^2$

- Заместваме и получаваме:

$$f(\theta_{k+1}) \leq f(\theta_k) - \alpha \|\nabla f(\theta_k)\|^2 + \frac{\alpha^2 L}{2} \|\nabla f(\theta_k)\|^2$$

- Нека подберем  $\alpha$ , така че  $\alpha L < 1$ . Тогава:

$$f(\theta_k) - f(\theta_{k+1}) \geq \frac{\alpha}{2} \|\nabla f(\theta_k)\|^2$$

- Следователно, на всяка стъпка **стойността на  $f$  намалява** и разликата на последователните членове на редицата е  $\geq \frac{\alpha}{2} \|\nabla f(\theta_k)\|^2$ .

- Нека сумираме нашето неравенство за  $k = 0, 1, \dots, T - 1$ :

$$\frac{\alpha}{2} \sum_{k=0}^{T-1} \|\nabla f(\theta_k)\|^2 \leq \sum_{k=0}^{T-1} f(\theta_k) - f(\theta_{k+1}) = f(\theta_0) - f(\theta_T)$$

- Нека  $f^*$  е глобален минимум на функцията  $f$ . Тогава

$$f(\theta_T) \geq f^* \text{ и следователно: } \sum_{k=0}^{T-1} \|\nabla f(\theta_k)\|^2 \leq \frac{2}{\alpha} (f(\theta_0) - f^*).$$

- Дясната страна е константа, следователно редът при  $T \rightarrow \infty$  с неотрицателни членове е ограничен, следователно е сходящ и  $\|\nabla f(\theta_k)\|^2 \rightarrow 0$ .
- $\|\theta_{k+1} - \theta_k\| = \alpha \|\nabla f(\theta_k)\| \rightarrow 0$ . Следователно редицата  $\theta_k$  (както и  $f(\theta_k)$ ) е сходяща. ■

# План на лекцията

---

1. Формалности за курса (5 мин)
2. Намиране на градиент чрез пропагиране назад — Backpropagation (20 мин)
3. Пропагиране назад при логистична регресия (20 мин)
4. Сходимость на спускането по градиента (20 мин)
5. **Стохастичен градиент (20 мин)**

# Стандартен стохастичен градиент

## Standard Stochastic Gradient

---

•  $f(\theta) = H_X[\text{Pr} \parallel \text{Pr}_\theta] = -\frac{1}{|X|} \sum_{i=1}^{|X|} \log \text{Pr}_\theta[\mathbf{x}^{(i)}]$ . Нека с

$H_{X_i}[\text{Pr} \parallel \text{Pr}_\theta] = -\log \text{Pr}_\theta[\mathbf{x}^{(i)}]$  означим поточковата кросентропия в точката  $X_i$  и  $f_{X_i}(\theta) = H_{X_i}[\text{Pr} \parallel \text{Pr}_\theta]$ .

• Нека  $X_{i_k}$  е случайно наблюдение (семпъл) от  $X$ . Тогава:

$$\theta_{k+1} = \theta_k - \alpha \nabla_{\theta} f_{X_{i_k}}(\theta_k)$$

• От теоремата на Тейлър получаваме:

$$f(\theta_{k+1}) = f(\theta_k) - \alpha \nabla f(\theta_k)^\top \nabla f_{X_{i_k}}(\theta_k) + \frac{1}{2} (\alpha \nabla f_{X_{i_k}}(\theta_k))^\top \nabla^2 f(\bar{\theta}) (\alpha \nabla f_{X_{i_k}}(\theta_k))$$

$$f(\theta_{k+1}) \leq f(\theta_k) - \alpha \nabla f(\theta_k)^\top \nabla f_{X_{i_k}}(\theta_k) + \frac{\alpha^2 L}{2} \|\nabla f_{X_{i_k}}(\theta_k)\|^2$$

# Свойства на стандартния стохастичен градиент

---

- Изчисляването на градиента става  $|X|$  пъти по-бързо.

- $$f(\theta_k) - f(\theta_{k+1}) \geq \alpha \nabla f(\theta_k)^\top \nabla f_{X_{i_k}}(\theta_k) - \frac{\alpha^2 L}{2} \|\nabla f_{X_{i_k}}(\theta_k)\|^2$$

- $\alpha \nabla f(\theta_k)^\top \nabla f_{X_{i_k}}(\theta_k)$  може да бъде и отрицателно, следователно нямаме никаква гаранция, че стойността на  $f$  намалява.

- Може да подходим **вероятностно**.

# Партиден стохастичен градиент

## Batched Stochastic Gradient

---

- Разглеждаме **партида** (batch, minibatch) — извадка от  $B$  на брой наблюдения на случайни величини от  $X$ :  $\mathbf{X}_B = X_1, X_2, \dots, X_B$ . Тогава ако означим с  $f_{\mathbf{X}_B}(\theta) = \frac{1}{B} \sum_{i=1}^B f_{X_i}(\theta_k)$  кросентропията на партидата  $\mathbf{X}_B$ , то дефинираме презаписването на параметрите като:
  - $\theta_{k+1} = \theta_k - \alpha \nabla_{\theta} f_{\mathbf{X}_B}(\theta_k)$ .
- Можем да повторим същите разсъждения като при стандартния стохастичен градиент в случая на партида. Разликите са, че времето за намиране на градиента нараства с фактор  $B$ , но за сметка на това може да очакваме, че отклонението на партидния градиента от ще намалее пълния градиент ще намалее.



# Заклучение

---

- Чрез Backpropagation градиентите на сложни функции се изчисляват автоматично, ефективно и точно, като в същото време се имплементират лесно.
- Спускането по пълния градиент е гарантирано сходящо, но се изчислява бавно.
- При стандартния стохастичен градиент не е гарантирано, че се намалява грешката, но е много по бързо.
- Партидният стохастичен градиент е компромис, при който има много по-висока вероятност за схождение, като е значително по-ефективен от пълния градиент.
- Партидният стохастичен градиент (и неговите вариации) е дефакто стандартният подход в съвременните системи за дълбоко машинно обучение.

# Логистична регресия — векторен запис

