

MAI4CAREU

Master programmes in Artificial
Intelligence 4 Careers in Europe



Πανεπιστήμιο Κύπρου - Τεχνητή Νοημοσύνη

MAI612 - ΜΗΧΑΝΙΚΗ ΜΑΘΗΣΗ

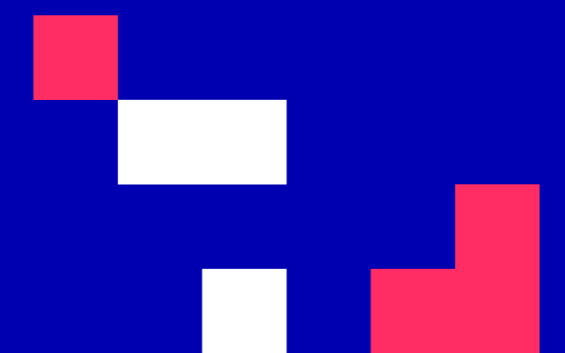
Διάλεξη 10: Νευρωνικά Δίκτυα 2: Εκπαίδευση

Βασίλης Βασιλειάδης, PhD

Χειμερινό Εξάμηνο 2022/23



CYENS
CENTRE OF EXCELLENCE





Διάλεξη 10: Νευρωνικά Δίκτυα 2: Εκπαίδευση

Μαθησιακά αποτελέσματα

Θα μάθετε για:

1. Η συνάρτηση κόστους των μοντέλων νευρωνικών δικτύων για παλινδρόμηση και ταξινόμηση
2. Πώς να εκπαιδεύσετε νευρωνικά δίκτυα χρησιμοποιώντας backpropagation
3. Πώς να εφαρμόσετε την οπισθοδρόμηση αποτελεσματικά
4. Στοχαστικό gradient descent με ορμή
5. Προηγμένες μέθοδοι βελτιστοποίησης για την εκπαίδευση των νευρωνικών δικτύων
6. Πώς να βελτιώσετε την απόδοση των νευρωνικών δικτύων χρησιμοποιώντας πρόωρη διακοπή, ρύθμιση υπερπαραμέτρου και σύνολα
7. Εξελισσόμενα νευρωνικά δίκτυα





Συνάρτηση κόστους παλινδρόμησης

Γραμμική παλινδρόμηση + Κανονικοποίηση:

$$L(\boldsymbol{\theta}) = \frac{1}{2m} \left[\sum_{i=1}^m (f_{\boldsymbol{\theta}}(x^{(i)}) - y^{(i)})^2 \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

Νευρωνικό δίκτυο + Κανονικοποίηση :

$$L(\Theta) = \frac{1}{2m} \left[\sum_{i=1}^m \sum_{k=1}^K \left((f_{\Theta}(x^{(i)}))_k - y_k^{(i)} \right)^2 \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{n_l} \sum_{j=1}^{n_{l+1}} \left(\Theta_{j,i}^{(l)} \right)^2$$

Κ μονάδες εξόδου
($\mathbf{y} \in \mathbb{R}^K$)





Συνάρτηση κόστους ταξινόμησης πολλαπλών κλάσεων

Πολυωνυμικό logistic regression+ Κανονισμός:

$$L(\boldsymbol{\theta}) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K \left[y_k^{(i)} \log(f_{\boldsymbol{\theta}}(x^{(i)}))_k \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

Νευρωνικό δίκτυο + Κανονισμός:

$$L(\Theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K \left[y_k^{(i)} \log(f_{\Theta}(x^{(i)}))_k \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{n_l} \sum_{j=1}^{n_{l+1}} \left(\Theta_{j,i}^{(l)} \right)^2$$





Υπολογισμός Gradient

Στόχος: ελαχιστοποιήστε $L(\Theta)$

Χρησιμοποιώντας **Gradient Descent**

Πρέπει να υπολογίσουμε:

$$L(\Theta)$$
$$\frac{\partial}{\partial \Theta_{j,i}^{(l)}} L(\Theta)$$

Παρατηρήσεις:

- Ένα ΝΔ έχει περισσότερες από 1 παράμετρο για να βελτιστοποιήσει
- Ένα ΝΔ είναι μια σύνθεση συναρτήσεων

Πρέπει να ξέρουμε πώς να υπολογίσουμε **μερικά παράγωγα** μιας **σύνθετης συνάρτησης**





Μερικά παράγωγα

- Για μια συνάρτηση δύο ή περισσότερων μεταβλητών, δεν $f(x_1, x_2)$ υπάρχει ένα μόνο παράγωγο, αλλά μερικά παράγωγα σε σχέση με κάθε μεταβλητή.
- Αυτά είναι μόνο τα συνήθη παράγωγα της συνάρτησης σε σχέση με μια μεταβλητή, διατηρώντας τις άλλες μεταβλητές σταθερές.
- Για παράδειγμα: για $f(x_1, x_2) = x_1^2 + 3x_2 - x_1x_2$

$$\frac{\partial f}{\partial x_1} = 2x_1 - x_2$$

$$\frac{\partial f}{\partial x_2} = 3 - x_1$$





Σύνθετες συναρτήσεις

- Για μια σύνθετη συνάρτηση $y = f[g(x)]$ μπορούμε να εφαρμόσουμε τον **κανόνα της αλυσίδας** για να πάρουμε $f'(x)$:

$$f'(x) = \frac{dy}{dx} = \frac{dy}{dg} \frac{dg}{dx} = f'[g(x)] g'(x)$$

- Αν και το επιχείρημα της f είναι μια άλλη $g(x)$ συνάρτηση, αντιμετωπίζουμε $g(x)$ σαν να ήταν απλά μια κανονική μεταβλητή και διαφοροποιείται f χρησιμοποιώντας τη συνήθη διαφοροποίηση σε σχέση με αυτή τη μεταβλητή.
- Στη συνέχεια πολλαπλασιάζουμε το παράγωγο του σε σχέση g με x , δηλαδή, $g'(x)$ για να πάρουμε το πλήρες παράγωγο.
- Ομοίως, αν $u = g(v)$ και $y = f(u)$ πού και $v = h(x)$ στη συνέχεια:

$$\frac{dy}{dx} = \frac{dy}{du} \frac{du}{dv} \frac{dv}{dx}$$





Σύνθετες λειτουργίες

Για παράδειγμα: $y = f[g(x)]$

$$g(x) = x^2$$
$$f(z) = 3z$$

$$z = g(x)$$
$$\Rightarrow f(g(x)) = f(x^2) = 3x^2$$
$$\Rightarrow f'(x) = \frac{df}{dx} = 6x$$

Κανόνας αλυσίδας:

$$\frac{dy}{dx} = \frac{dy}{dg} \frac{dg}{dx}$$
$$= 3 \cdot 2x$$
$$= 6x$$





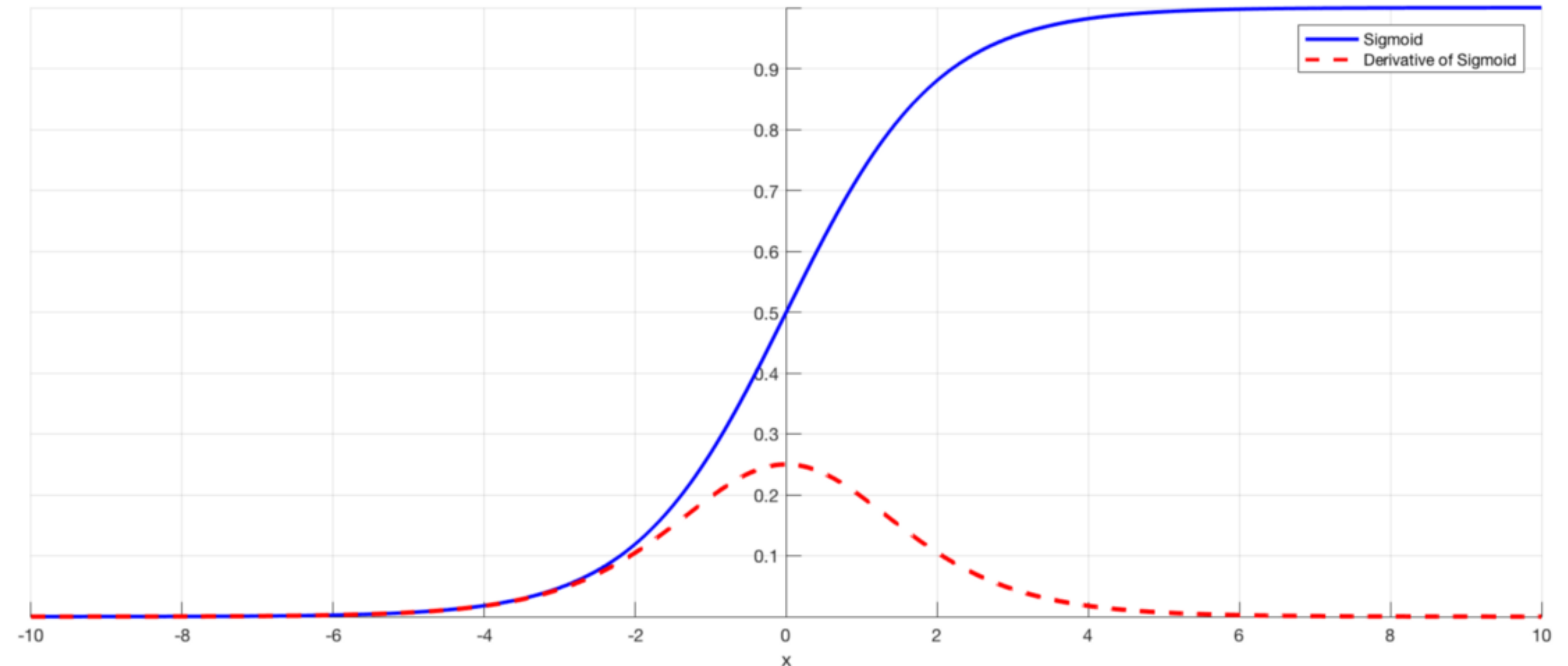
Παράγωγος της σιγμοειδούς συνάρτησης

$$y = f(x) = \frac{1}{1 + \exp(-ax)}$$

$$f'(x) = \frac{d}{dx} \left(\frac{1}{1 + \exp(-ax)} \right)$$

$$= \frac{a \exp(-ax)}{(1 + \exp(-ax))^2}$$

$$= a f(x) (1 - f(x))$$





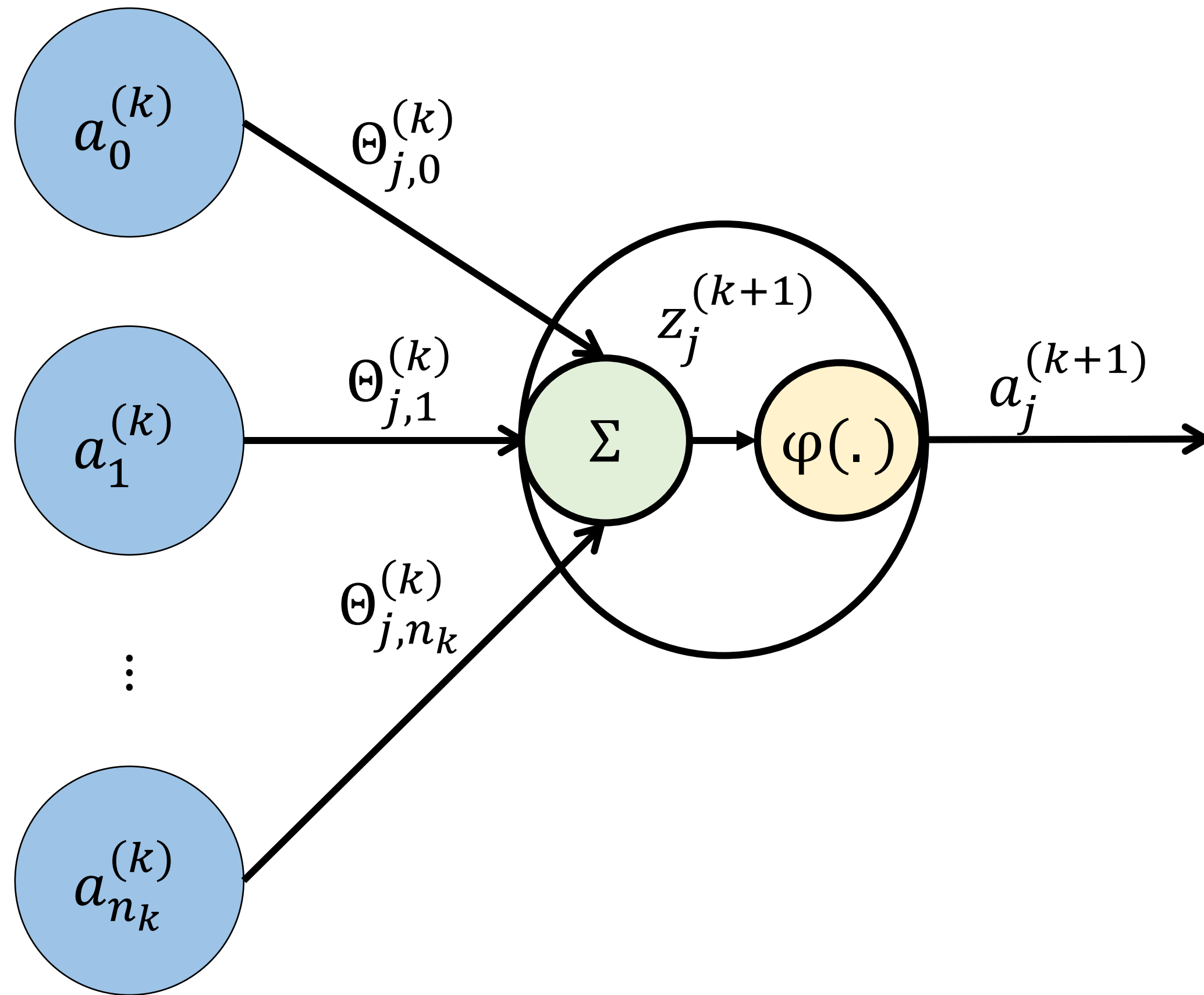
Αλγόριθμος οπισθοδιάδοσης

- **Διάδοση προς τα εμπρός:**
 - Παίρνει ένα μοτίβο εισόδου και το μετατρέπει διαδίδοντάς το προς τα εμπρός στο δίκτυο
 - Υπολογίζει την ενεργοποίηση κάθε κόμβου $a_i^{(k)}$ για να υπολογίσει το σφάλμα
- **Οπισθοδιάδοση:**
 - Υπολογίζει την κλίση σε ένα συγκεκριμένο μοτίβο εισόδου
 - «Αντιστρέφει τη ροή των πληροφοριών» διαδίδοντας τα **σφάλματα** προς τα πίσω στο δίκτυο (χρησιμοποιεί τον κανόνα της αλυσίδας)
 - Διαισθητικά: κάθε κόμβος πρέπει να υπολογίσει το δικό του σφάλμα: $\delta_i^{(k)}$
- Η οπίσθια διάδοση αναφέρεται μερικές φορές ως αλγόριθμος που εκπαιδεύει ένα νευρωνικό δίκτυο.
 - Είναι στην πραγματικότητα **ένας αποτελεσματικός** τρόπος υπολογισμού της κλίσης
 - Μπορούν να χρησιμοποιηθούν διάφορες μέθοδοι βελτιστοποίησης βάσει κλίσεων (συνήθως κάθοδος κλίσης)





Διάδοση προς τα εμπρός



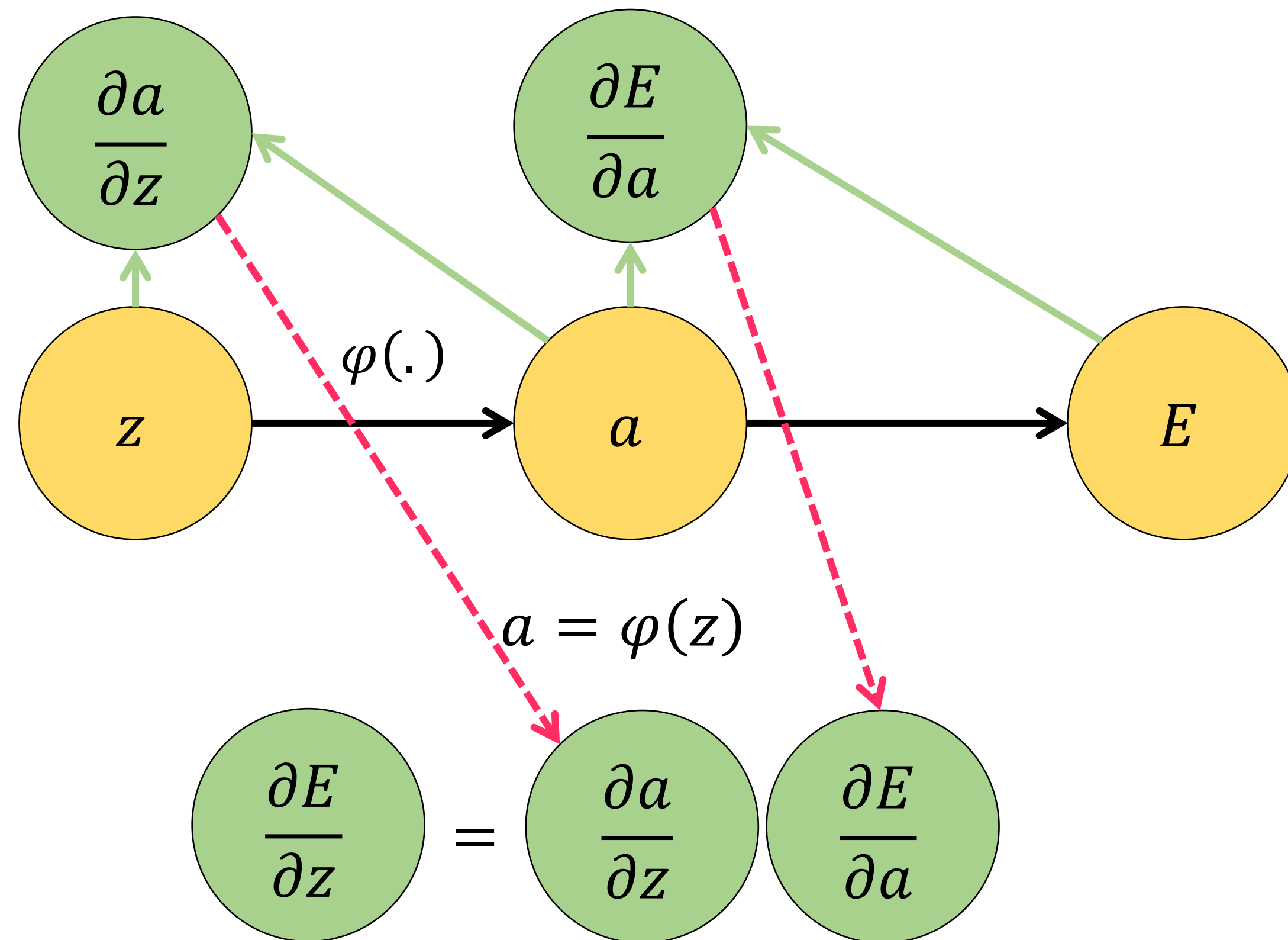
$$z_j^{(k+1)} = \sum_{i=0}^{n_k} \Theta_{j,i}^{(k)} a_i^{(k)}$$

$$a_j^{(k+1)} = \varphi \left(z_j^{(k+1)} \right)$$





Οπισθοδιάδοση



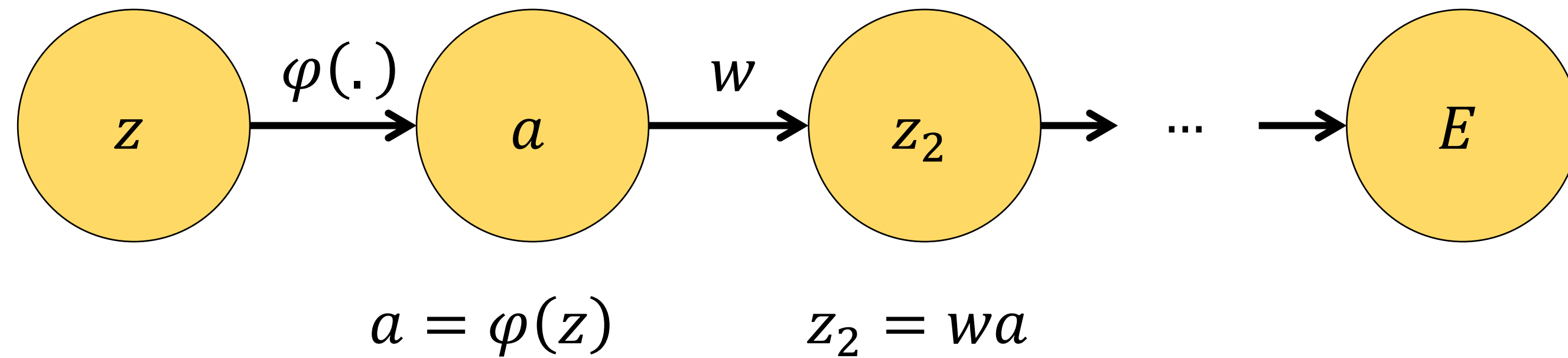
$$\frac{\partial E}{\partial z} = \frac{\partial E}{\partial a} \frac{\partial a}{\partial z}$$

$$\frac{\partial E}{\partial z} = \frac{\partial E}{\partial a} \varphi'(z)$$





Οπισθοδιάδοση



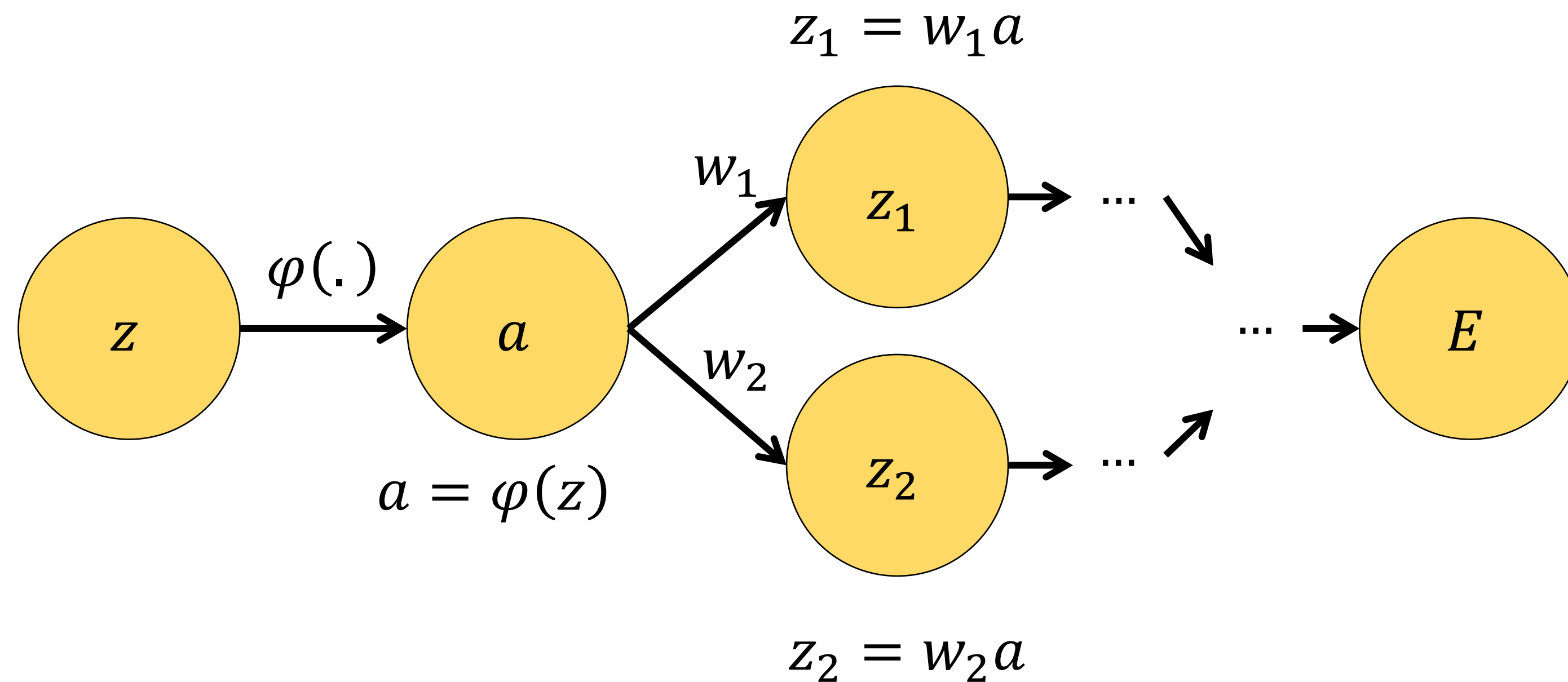
$$\frac{\partial E}{\partial z} = \frac{\partial E}{\partial z_2} \frac{\partial z_2}{\partial a} \frac{\partial a}{\partial z}$$

$$\frac{\partial E}{\partial z} = \frac{\partial E}{\partial z_2} w \varphi'(z)$$





Οπισθοδιάδοση



$$\frac{\partial E}{\partial z} = \left(\frac{\partial E}{\partial z_1} \frac{\partial z_1}{\partial a} + \frac{\partial E}{\partial z_2} \frac{\partial z_2}{\partial a} \right) \frac{\partial a}{\partial z}$$

$$\frac{\partial E}{\partial z} = \varphi'(z) \sum_{k=1}^{\#out} w_k \frac{\partial E}{\partial z_k}$$





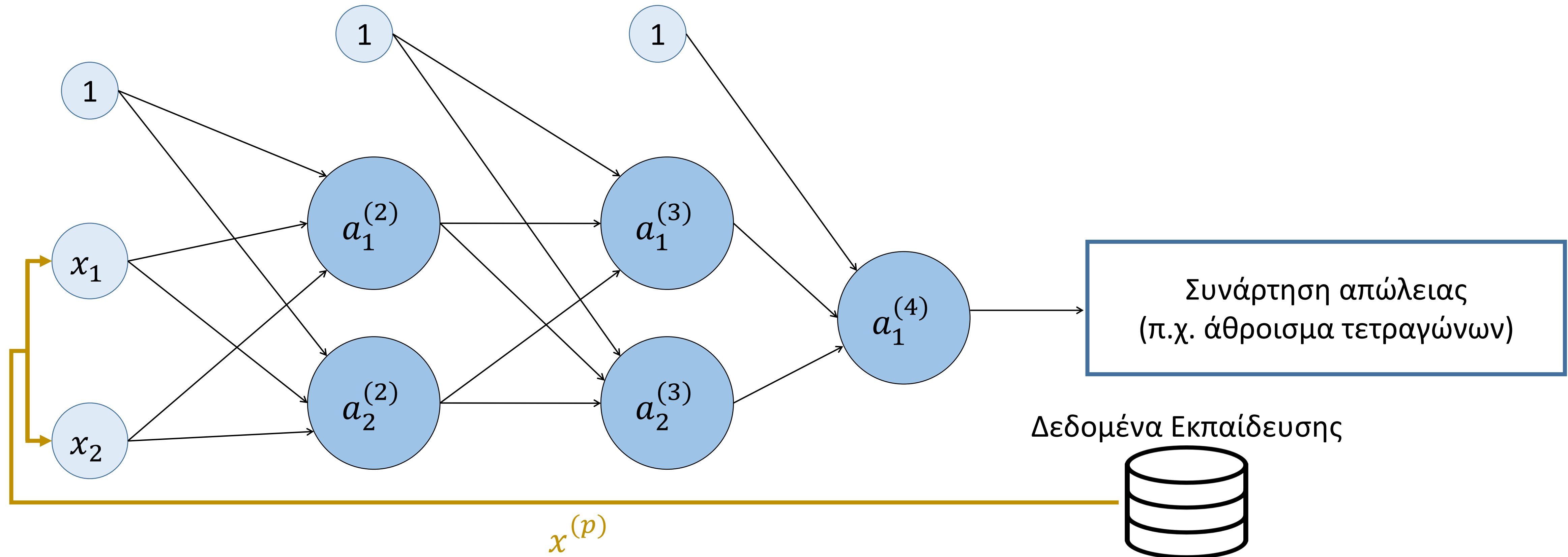
Αλγόριθμος οπισθοδιάδοσης Online Updating

Βήμα προς βήμα



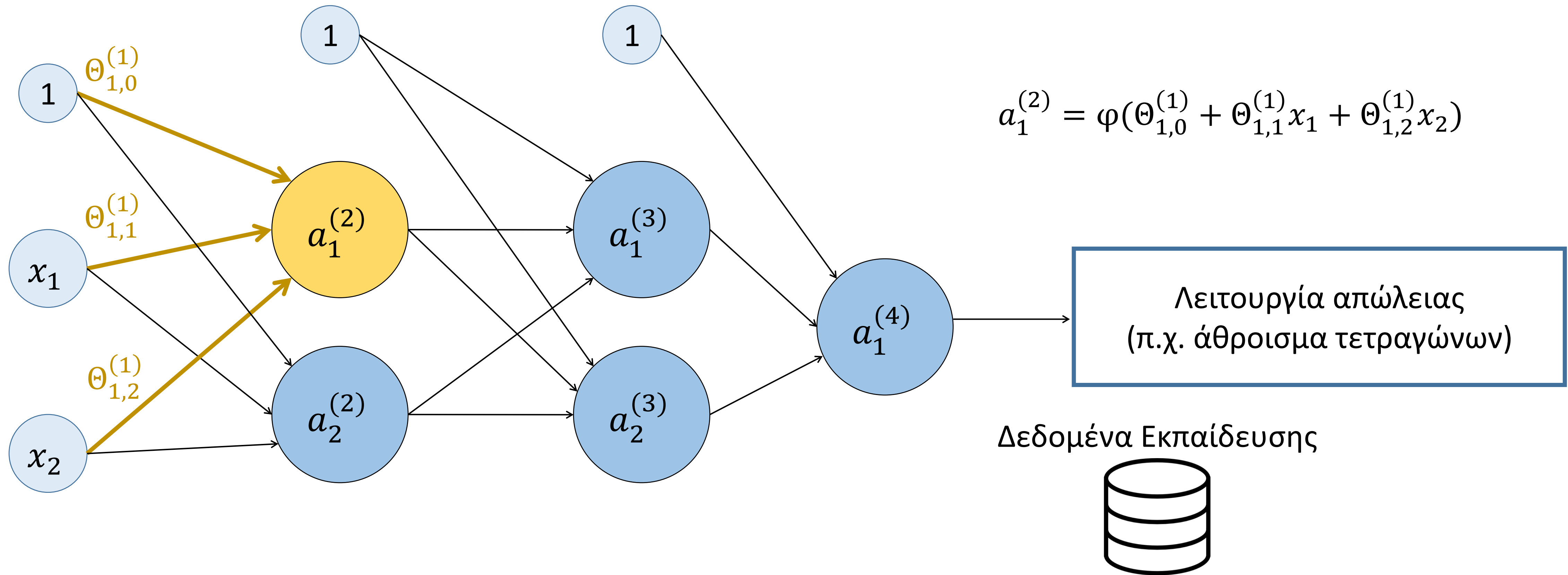


Το μοντέλο p παρουσιάζεται



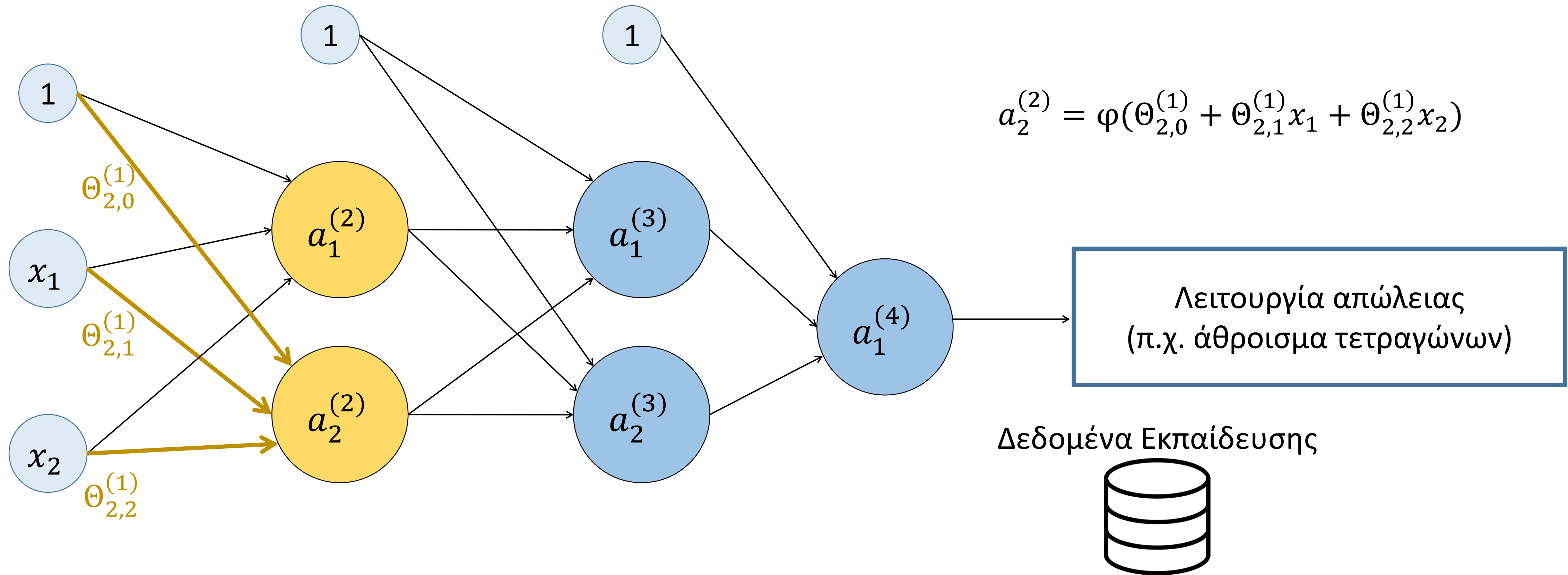


Διάδοση προς τα εμπρός: Υπολογισμός κόμβων εξόδου



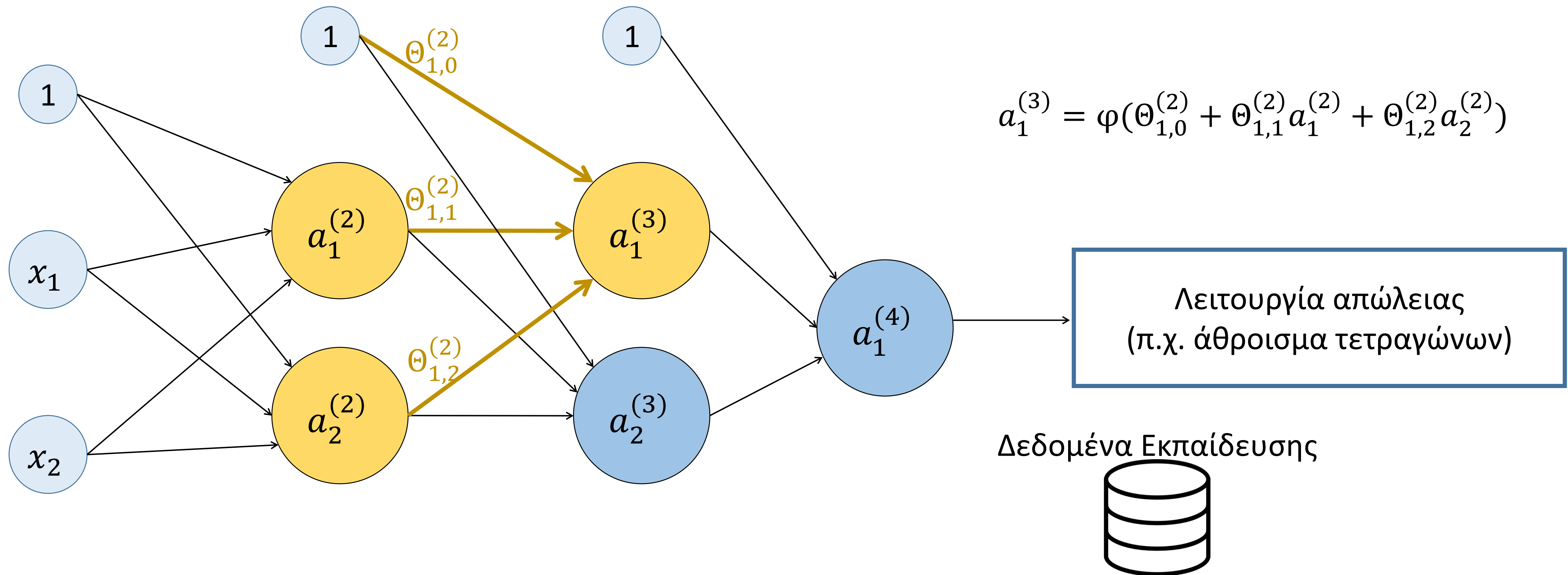


Διάδοση προς τα εμπρός: Υπολογισμός κόμβων εξόδου





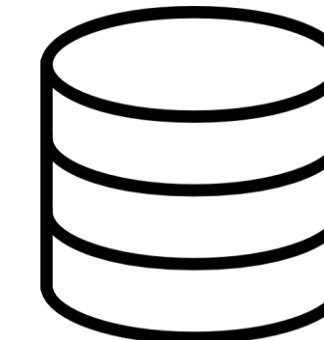
Διάδοση προς τα εμπρός: Υπολογισμός κόμβων εξόδου



$$a_1^{(3)} = \varphi(\Theta_{1,0}^{(2)} + \Theta_{1,1}^{(2)} a_1^{(2)} + \Theta_{1,2}^{(2)} a_2^{(2)})$$

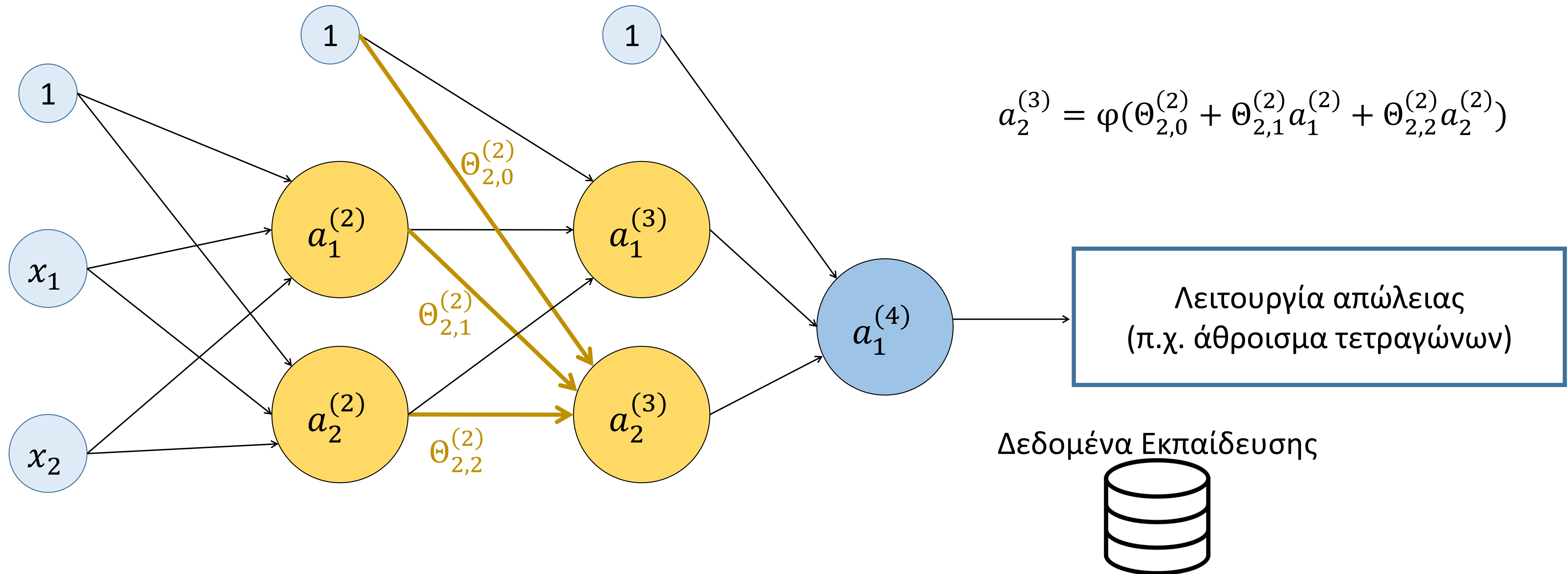
Λειτουργία απώλειας
(π.χ. άθροισμα τετραγώνων)

Δεδομένα Εκπαίδευσης



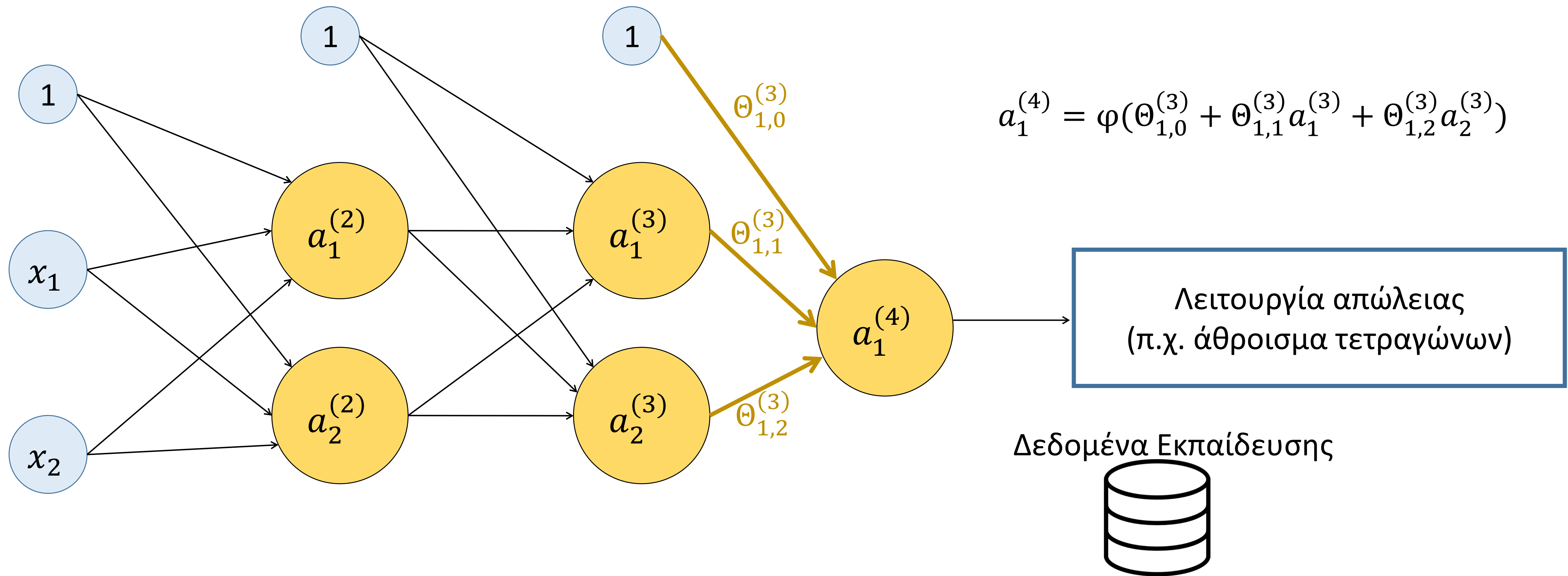


Διάδοση προς τα εμπρός: Υπολογισμός κόμβων εξόδου





Διάδοση προς τα εμπρός: Υπολογισμός κόμβων εξόδου



$$a_1^{(4)} = \varphi(\Theta_{1,0}^{(3)} + \Theta_{1,1}^{(3)} a_1^{(3)} + \Theta_{1,2}^{(3)} a_2^{(3)})$$

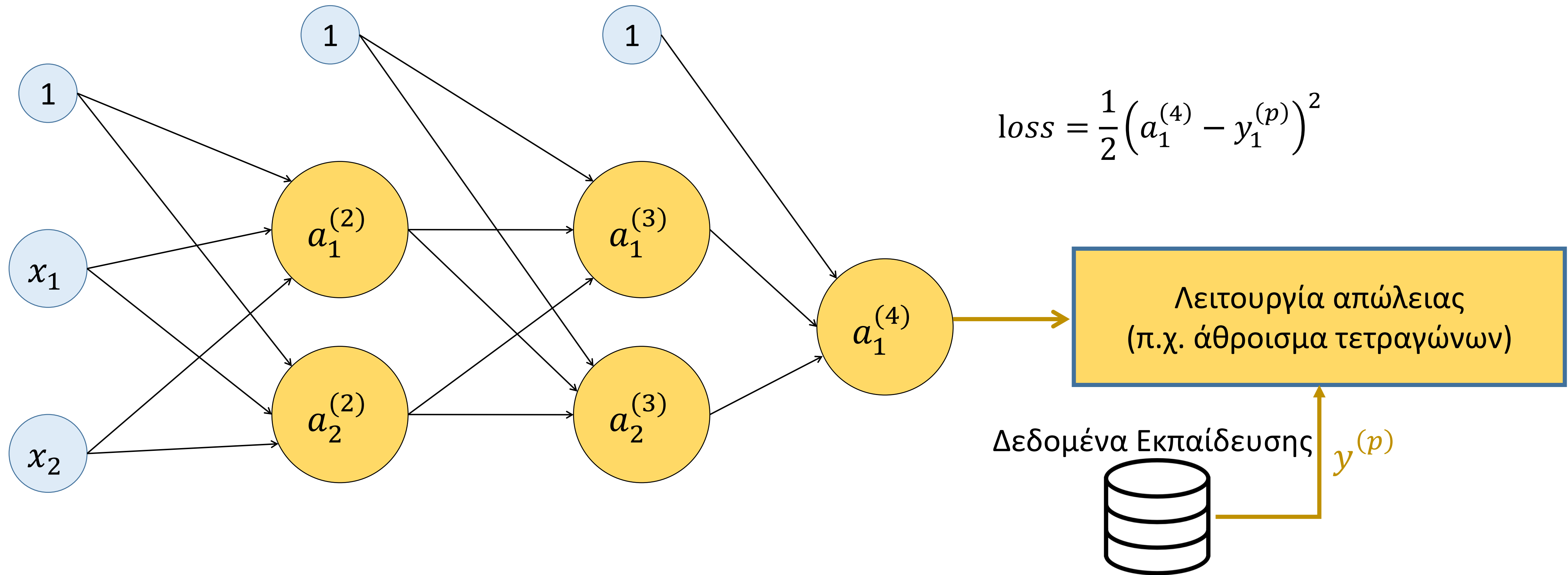
Λειτουργία απώλειας
(π.χ. άθροισμα τετραγώνων)

Δεδομένα Εκπαίδευσης



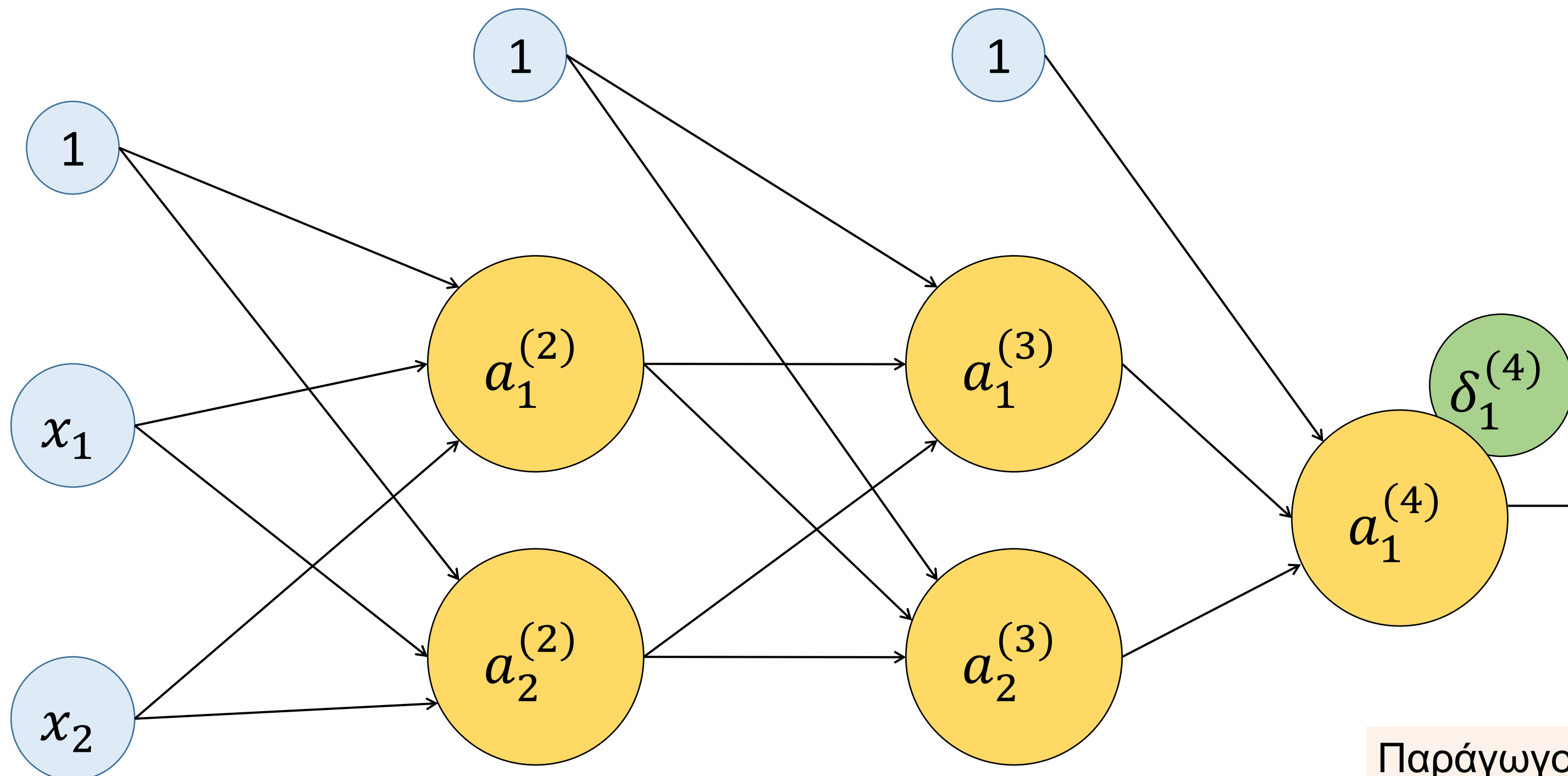


Διάδοση προς τα εμπρός: Υπολογισμός κόμβων εξόδου





Οπισθοδιάδοση: Υπολογισμός μερικών παραγώγων



Παράγωγο της ζημίας σε σχέση με την παραγωγή $a_1^{(4)}$

$$\delta_1^{(4)} = a_1^{(4)} (1 - a_1^{(4)}) \cdot (a_1^{(4)} - y_1^{(p)})$$

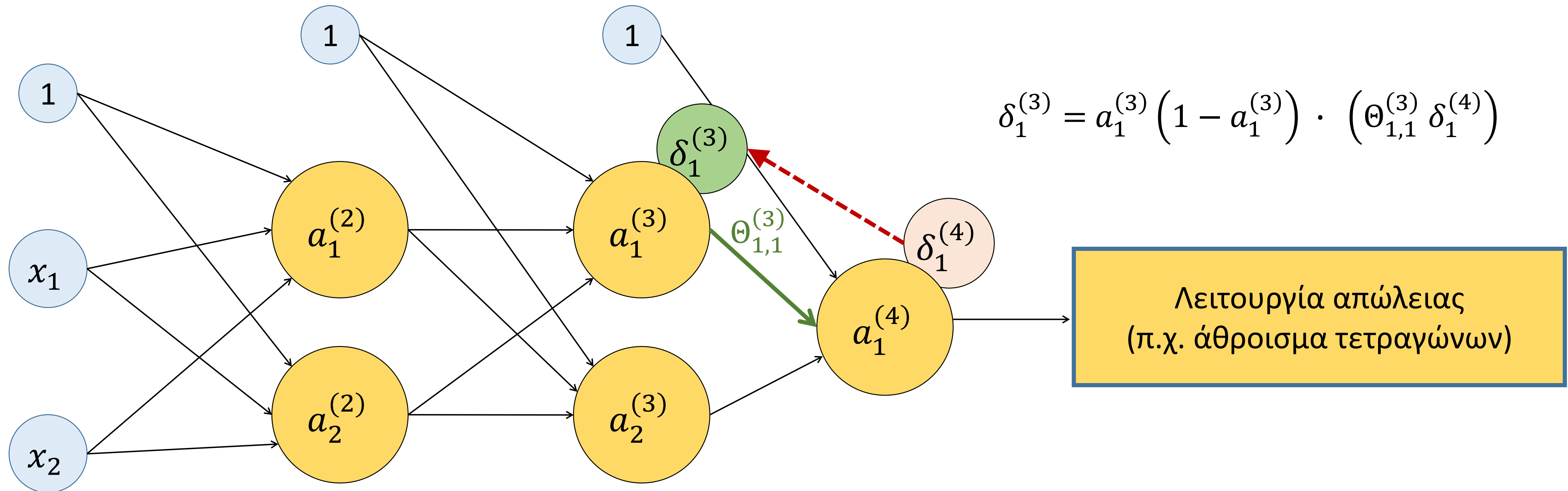
Λειτουργία απώλειας (π.χ. άθροισμα τετραγώνων)

Παράγωγος της λειτουργίας ενεργοποίησης του κόμβου εξόδου (δηλαδή, σιγμοειδής)





Οπισθοδιάδοση: Υπολογισμός μερικών παραγώγων



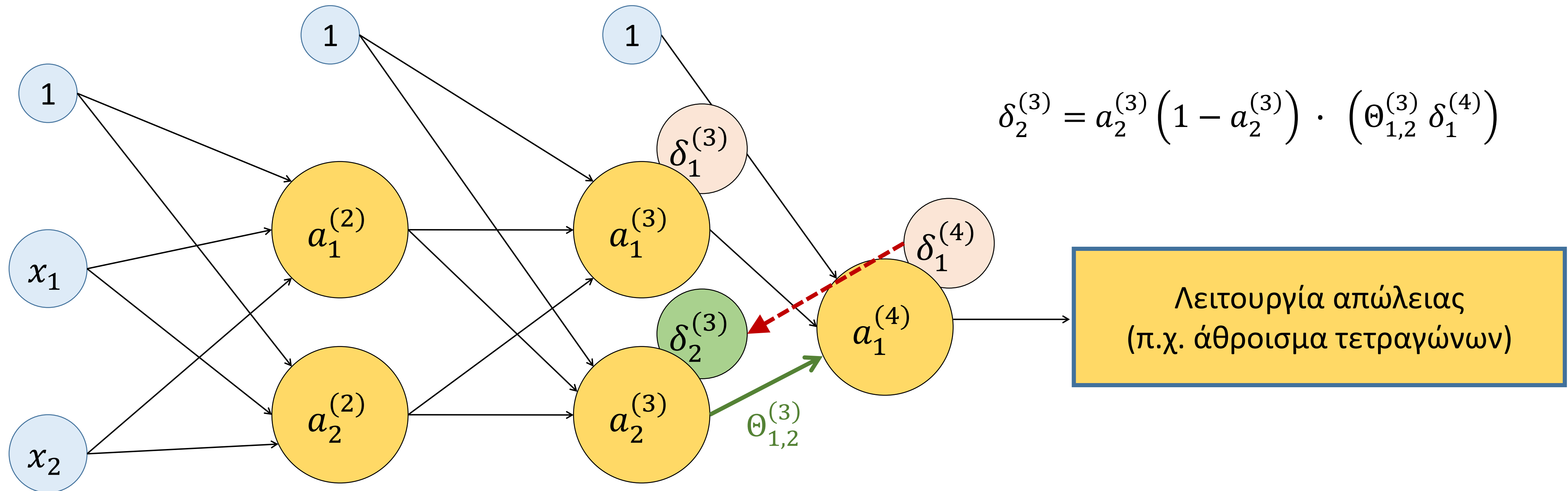
$$\delta_1^{(3)} = a_1^{(3)} (1 - a_1^{(3)}) \cdot (\Theta_{1,1}^{(3)} \delta_1^{(4)})$$

Λειτουργία απώλειας
(π.χ. άθροισμα τετραγώνων)





Οπισθοδιάδοση: Υπολογισμός μερικών παραγώγων



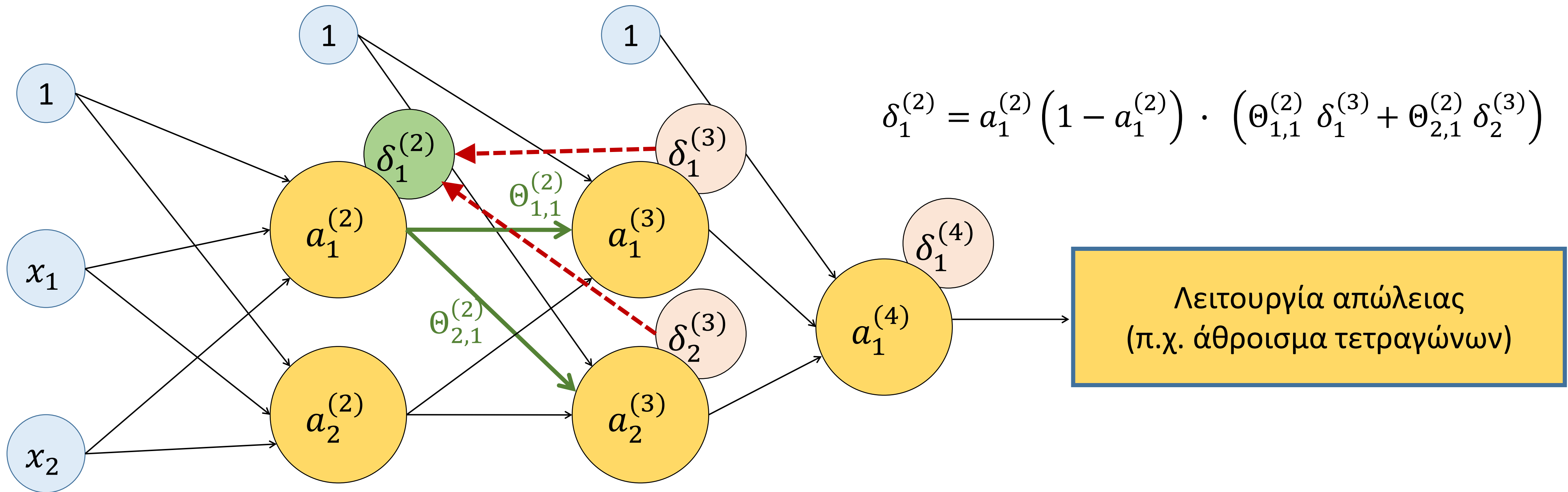
$$\delta_2^{(3)} = a_2^{(3)} (1 - a_2^{(3)}) \cdot (\Theta_{1,2}^{(3)} \delta_1^{(4)})$$

Λειτουργία απώλειας
(π.χ. άθροισμα τετραγώνων)





Οπισθοδιάδοση: Υπολογισμός μερικών παραγώγων



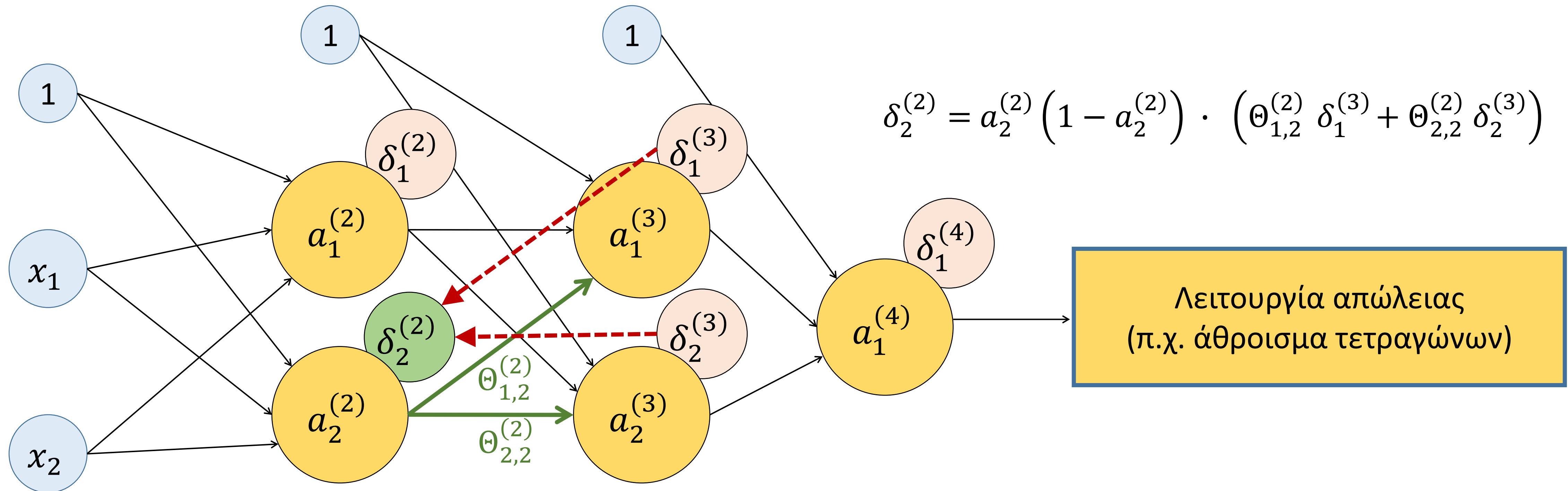
$$\delta_1^{(2)} = a_1^{(2)} (1 - a_1^{(2)}) \cdot (\Theta_{1,1}^{(2)} \delta_1^{(3)} + \Theta_{2,1}^{(2)} \delta_2^{(3)})$$

Λειτουργία απώλειας
(π.χ. άθροισμα τετραγώνων)



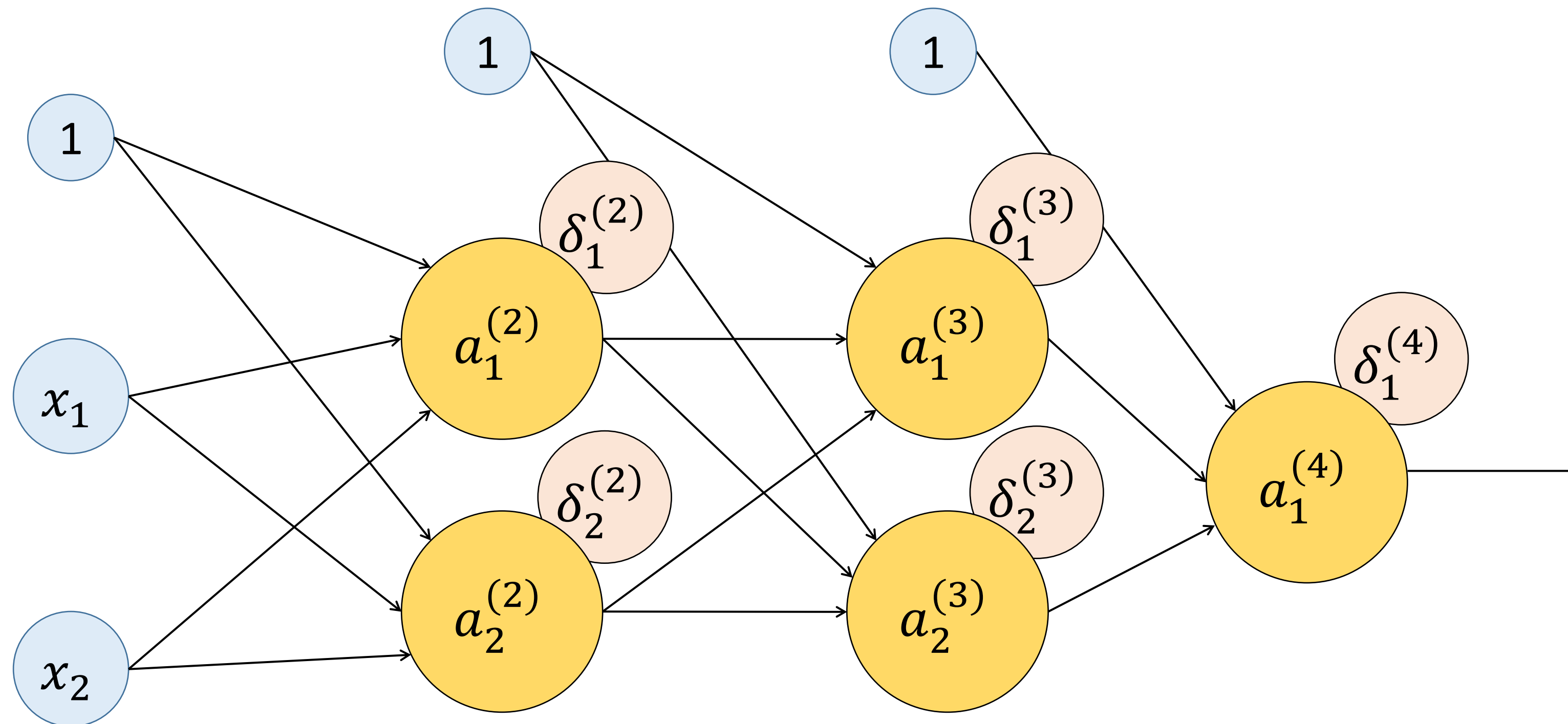


Οπισθοδιάδοση: Υπολογισμός μερικών παραγώγων



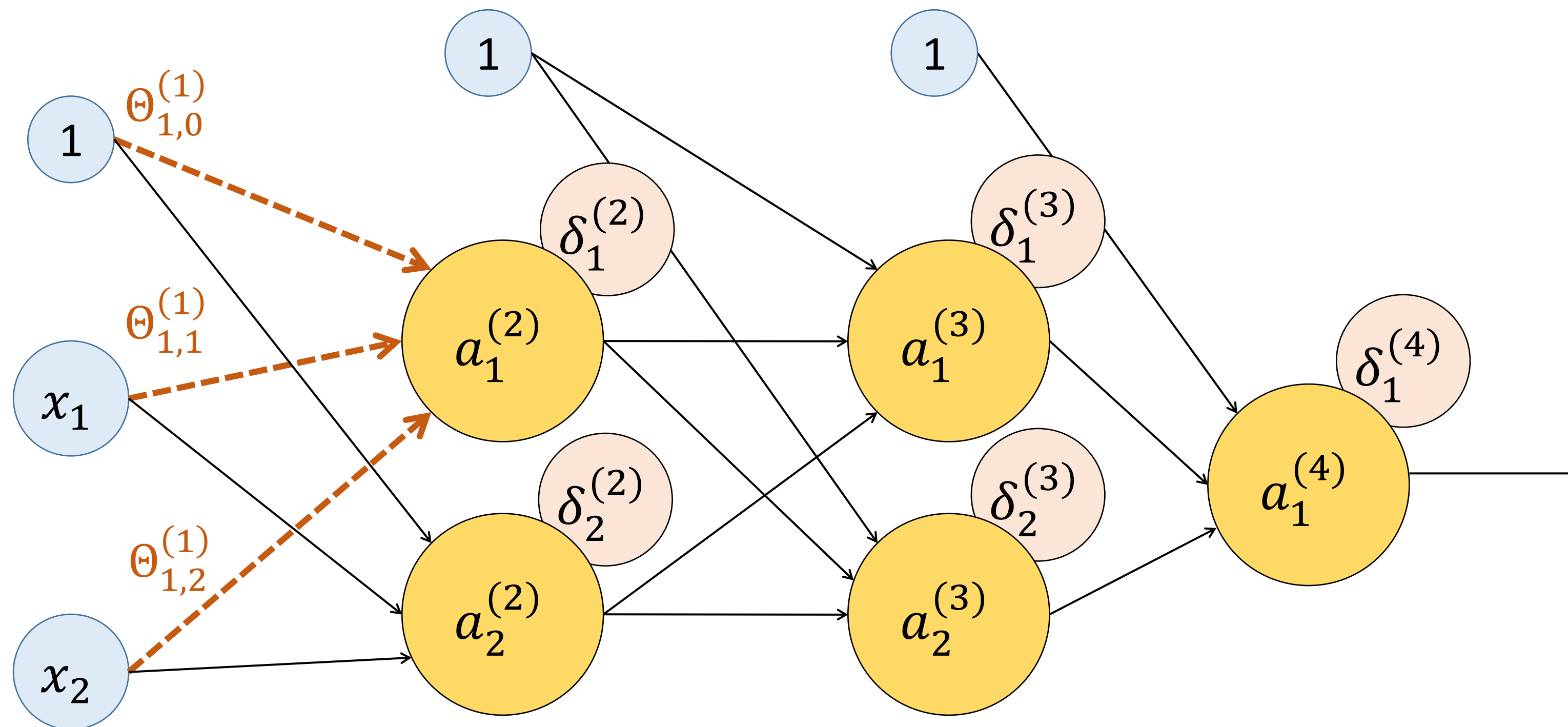


Ενημέρωση των παραμέτρων με στοχαστικό gradient descent





Ενημέρωση των παραμέτρων με στοχαστικό gradient descent



Παράγωγη της απώλειας σε σχέση με $\Theta_{1,0}^{(1)}$

$$\Theta_{1,0}^{(1)} := \Theta_{1,0}^{(1)} - \eta \delta_1^{(2)} x_0$$

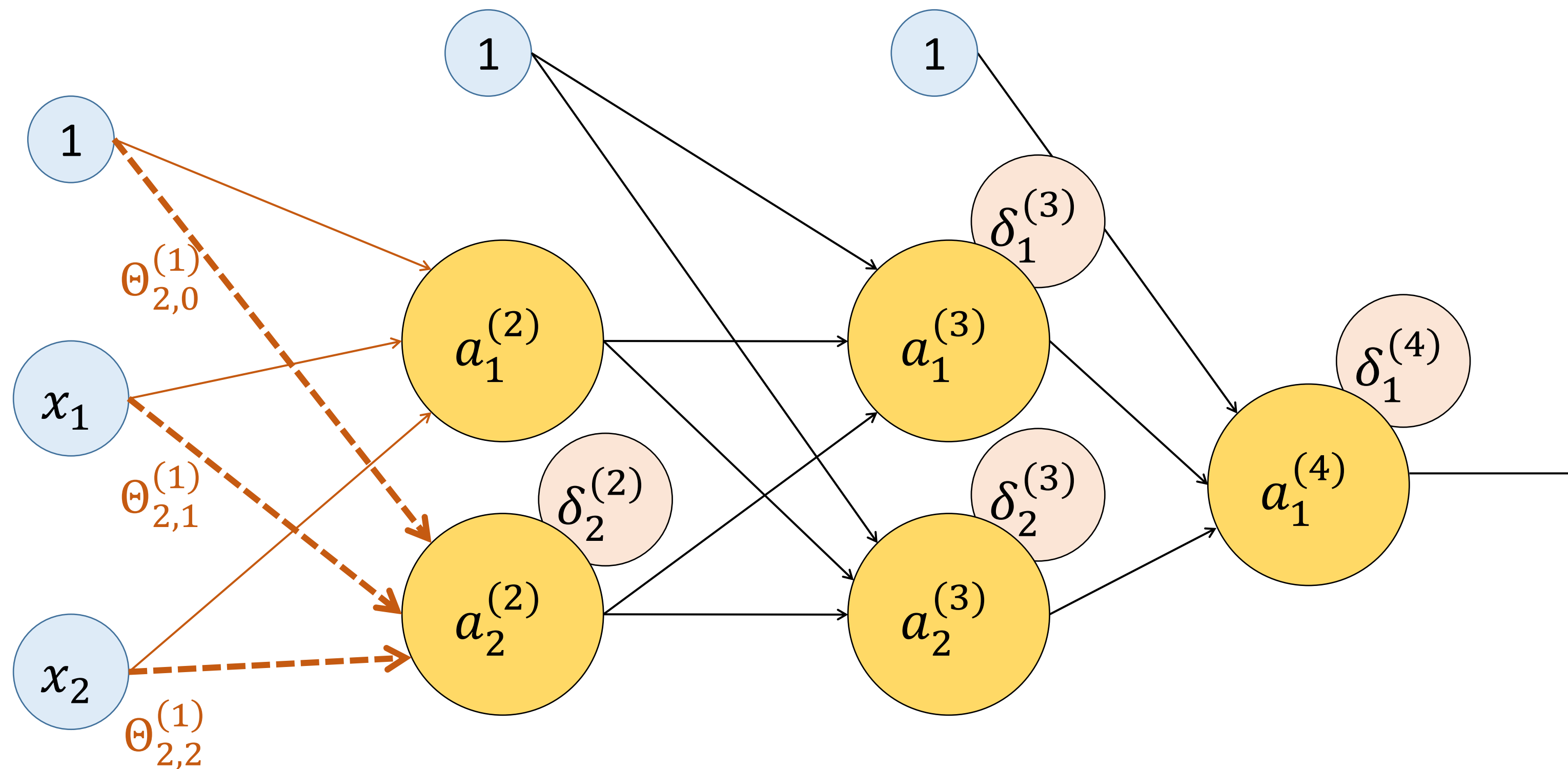
$$\Theta_{1,1}^{(1)} := \Theta_{1,1}^{(1)} - \eta \delta_1^{(2)} x_1$$

$$\Theta_{1,2}^{(1)} := \Theta_{1,2}^{(1)} - \eta \delta_1^{(2)} x_2$$





Ενημέρωση των παραμέτρων με στοχαστικό gradient descent



$$\Theta_{2,0}^{(1)} := \Theta_{2,0}^{(1)} - \eta \delta_2^{(2)} x_0$$

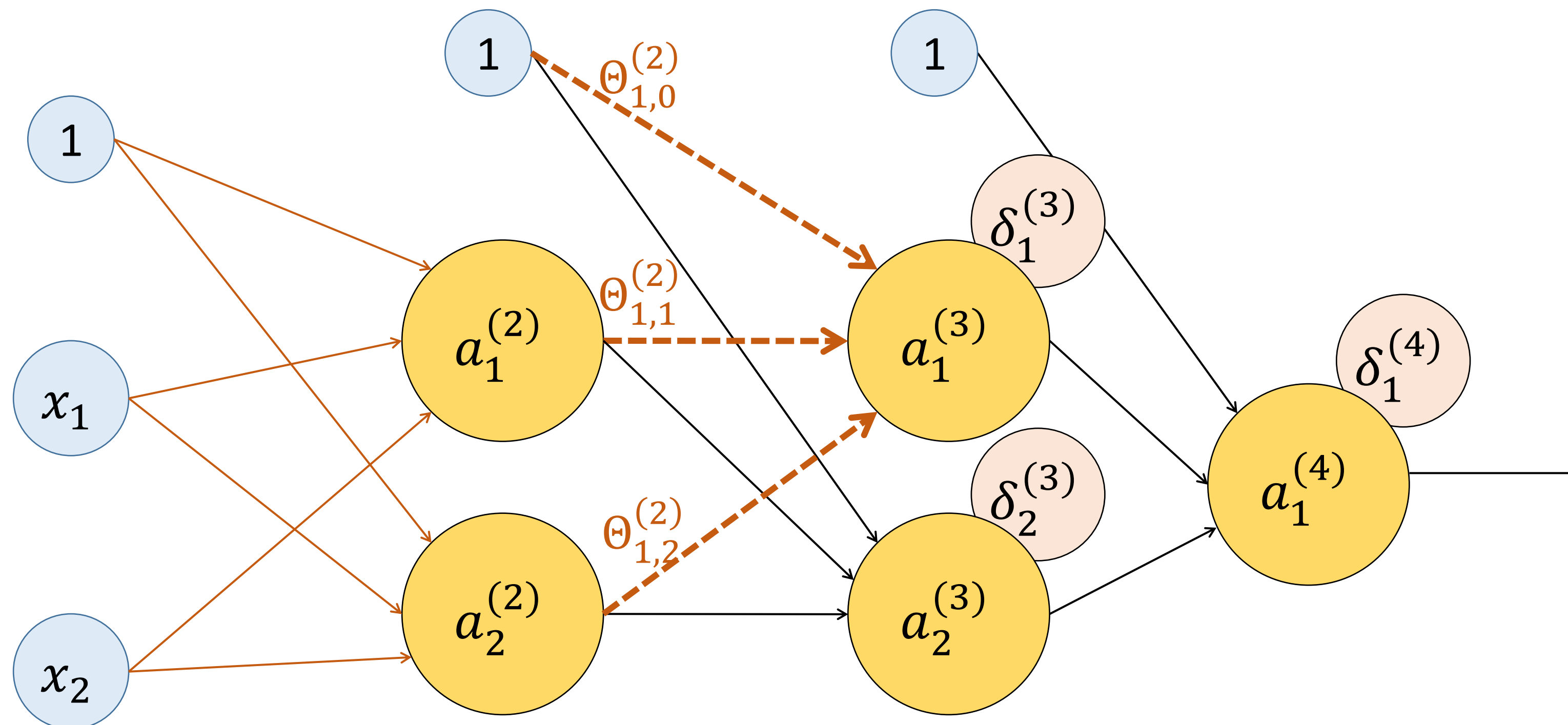
$$\Theta_{2,1}^{(1)} := \Theta_{2,1}^{(1)} - \eta \delta_2^{(2)} x_1$$

$$\Theta_{2,2}^{(1)} := \Theta_{2,2}^{(1)} - \eta \delta_2^{(2)} x_2$$





Ενημέρωση των παραμέτρων με στοχαστικό gradient descent



$$\Theta_{1,0}^{(2)} := \Theta_{1,0}^{(2)} - \eta \delta_1^{(3)} a_0^{(2)}$$

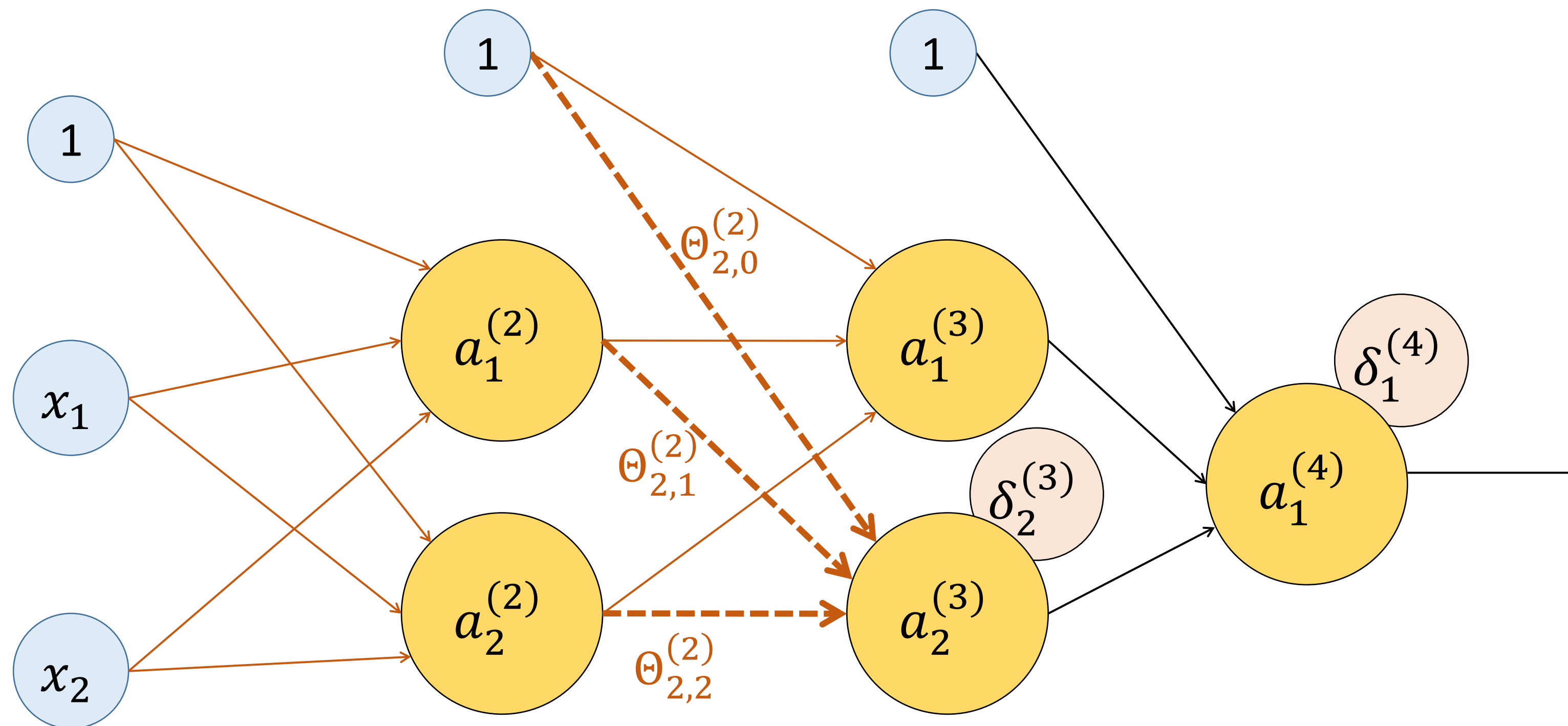
$$\Theta_{1,1}^{(2)} := \Theta_{1,1}^{(2)} - \eta \delta_1^{(3)} a_1^{(2)}$$

$$\Theta_{1,2}^{(2)} := \Theta_{1,2}^{(2)} - \eta \delta_1^{(3)} a_2^{(2)}$$





Ενημέρωση των παραμέτρων με στοχαστικό gradient descent



$$\Theta_{2,0}^{(2)} := \Theta_{2,0}^{(2)} - \eta \delta_2^{(3)} a_0^{(2)}$$

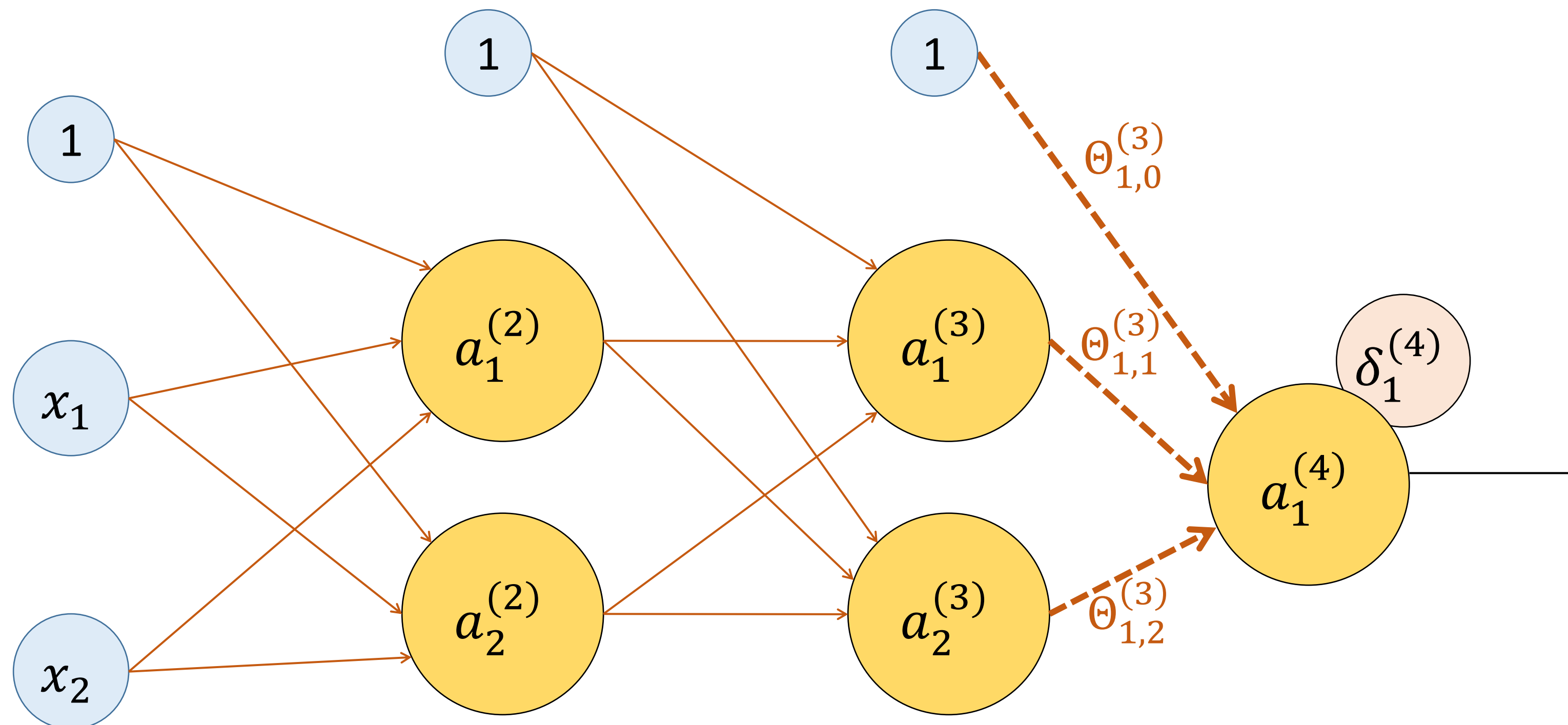
$$\Theta_{2,1}^{(2)} := \Theta_{2,1}^{(2)} - \eta \delta_2^{(3)} a_1^{(2)}$$

$$\Theta_{2,2}^{(2)} := \Theta_{2,2}^{(2)} - \eta \delta_2^{(3)} a_2^{(2)}$$





Ενημέρωση των παραμέτρων με στοχαστικό gradient descent



$$\Theta_{1,0}^{(3)} := \Theta_{1,0}^{(3)} - \eta \delta_1^{(4)} a_0^{(3)}$$

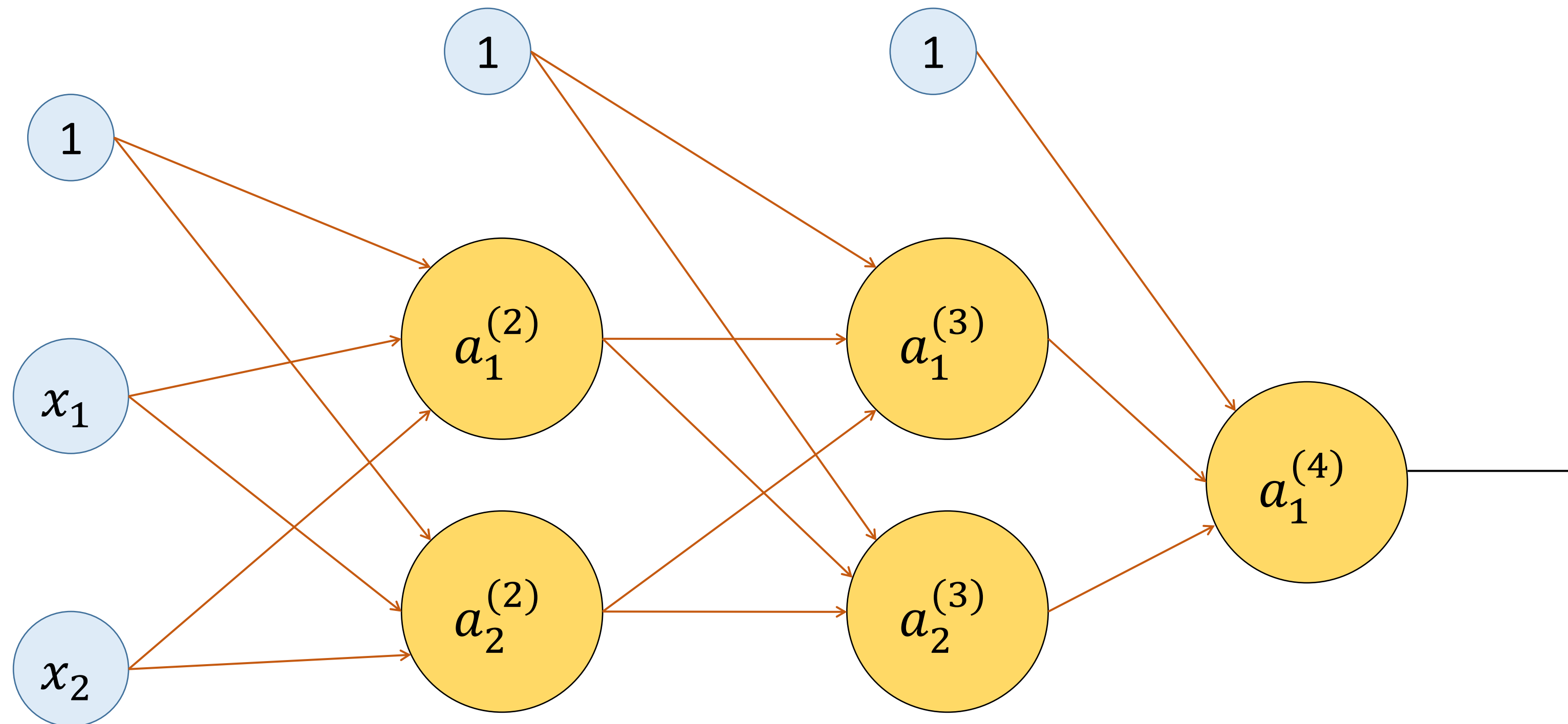
$$\Theta_{1,1}^{(3)} := \Theta_{1,1}^{(3)} - \eta \delta_1^{(4)} a_1^{(3)}$$

$$\Theta_{1,2}^{(3)} := \Theta_{1,2}^{(3)} - \eta \delta_1^{(4)} a_2^{(3)}$$



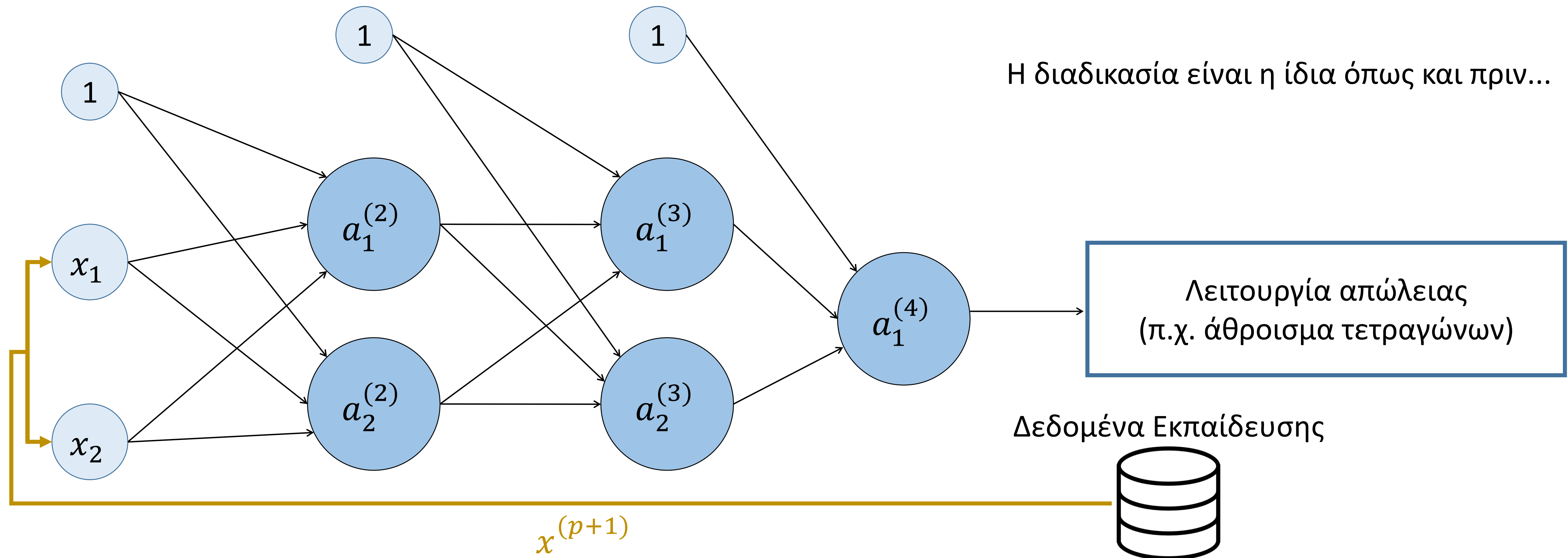


Όλες οι παράμετροι έχουν επικαιροποιηθεί





Παρουσιάζεται το μοτίβο $p+1$





Αλγόριθμος οπισθοδιάδοσης

Εξίσωση επικαιροποίησης βάρους:

$$\Delta \Theta_{j,i}^{(k)} = -\eta \cdot \frac{\partial}{\partial \Theta_{j,i}^{(l)}} L(\Theta)$$

Online updating
(stochastic gradient descent)

ΠΟΥ:

$$\frac{\partial}{\partial \Theta_{j,i}^{(l)}} L(\Theta) = \delta_j^{(k+1)} a_i^{(k)}$$

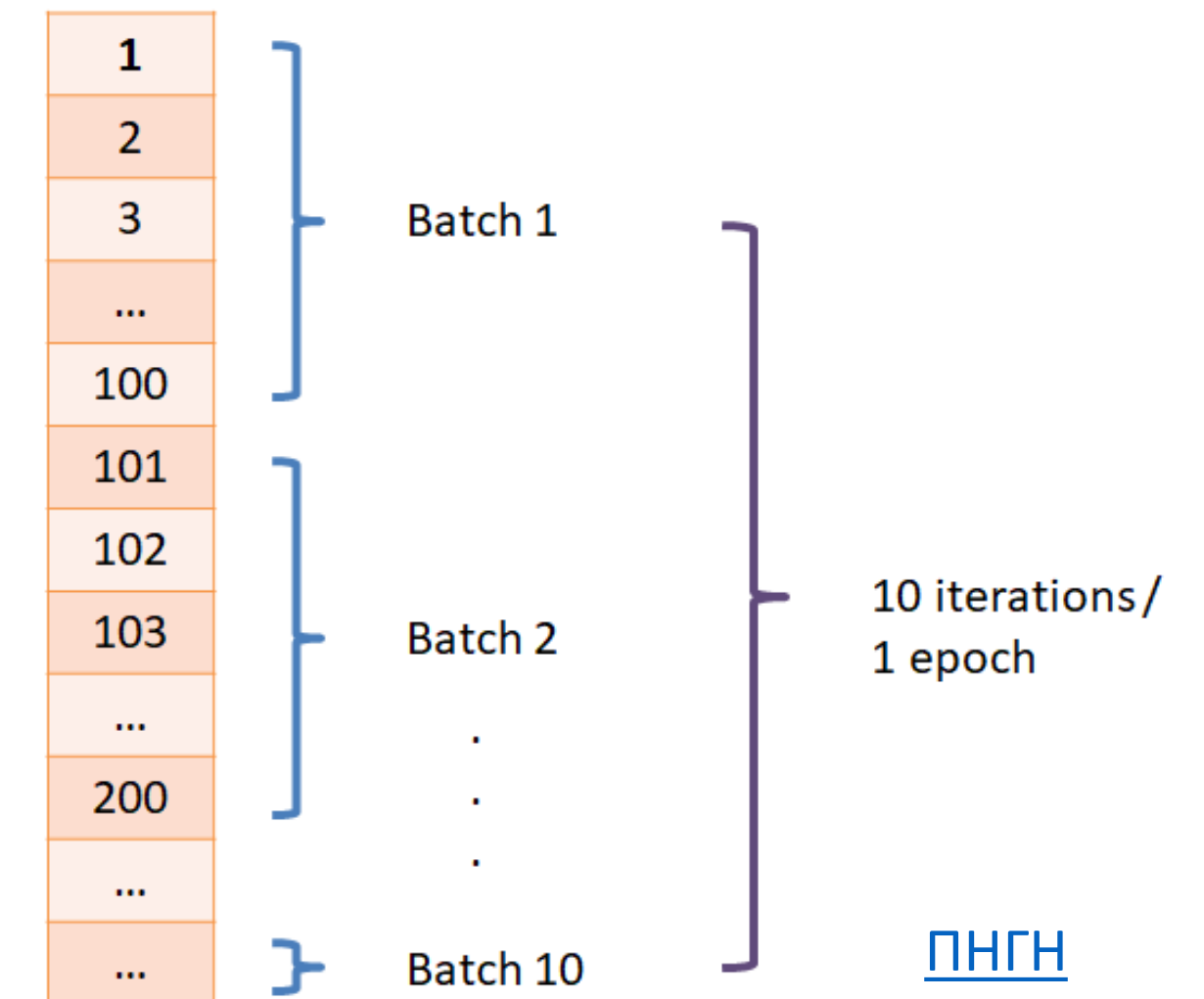
$$\frac{\partial}{\partial \Theta_{j,i}^{(l)}} L(\Theta) = \frac{1}{m} \sum_{p=1}^m \delta_j^{(k+1)}(p) a_i^{(k)}(p)$$

Batch updating
(batch gradient descent)

Δειγματοληψία/επικαιροποίηση mini-batch:

- Αντί να χρησιμοποιείτε ολόκληρο το σύνολο εκπαίδευσης, χρησιμοποιήστε μη επικαλυπτόμενα υποσύνολα του (mini-batches)
- Χρησιμοποιήστε τη διαδικασία επικαιροποίησης της παρτίδας m , όπου m είναι το μέγεθος του mini-batch

All training samples





Τροφοδοσία κατάρτισης ΝΔ με χρήση οπισθοδιάδοσης

Για κάθε εποχή (ή μέχρι το συνολικό σφάλμα ή το gradient norm να είναι αρκετά μικρό)

Θέσε $\Delta_{ij}^{(l)} = 0$ (για όλους l, i, j)

για $i = 1$ μέχρι m

Θέσε $a^{(1)} = x^{(i)}$

Εμπρόσθια διάδοση: Υπολογισμός $a^{(l)}$ για $l = 2, 3, \dots, L$

Οπίσθια διάδοση: Υπολογισμός $\delta^{(L)} = a^{(L)} - y^{(i)}, \delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$

$$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$$

$$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)} \quad \text{εάν } j \neq 0$$

$$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} \quad \text{εάν } j = 0$$

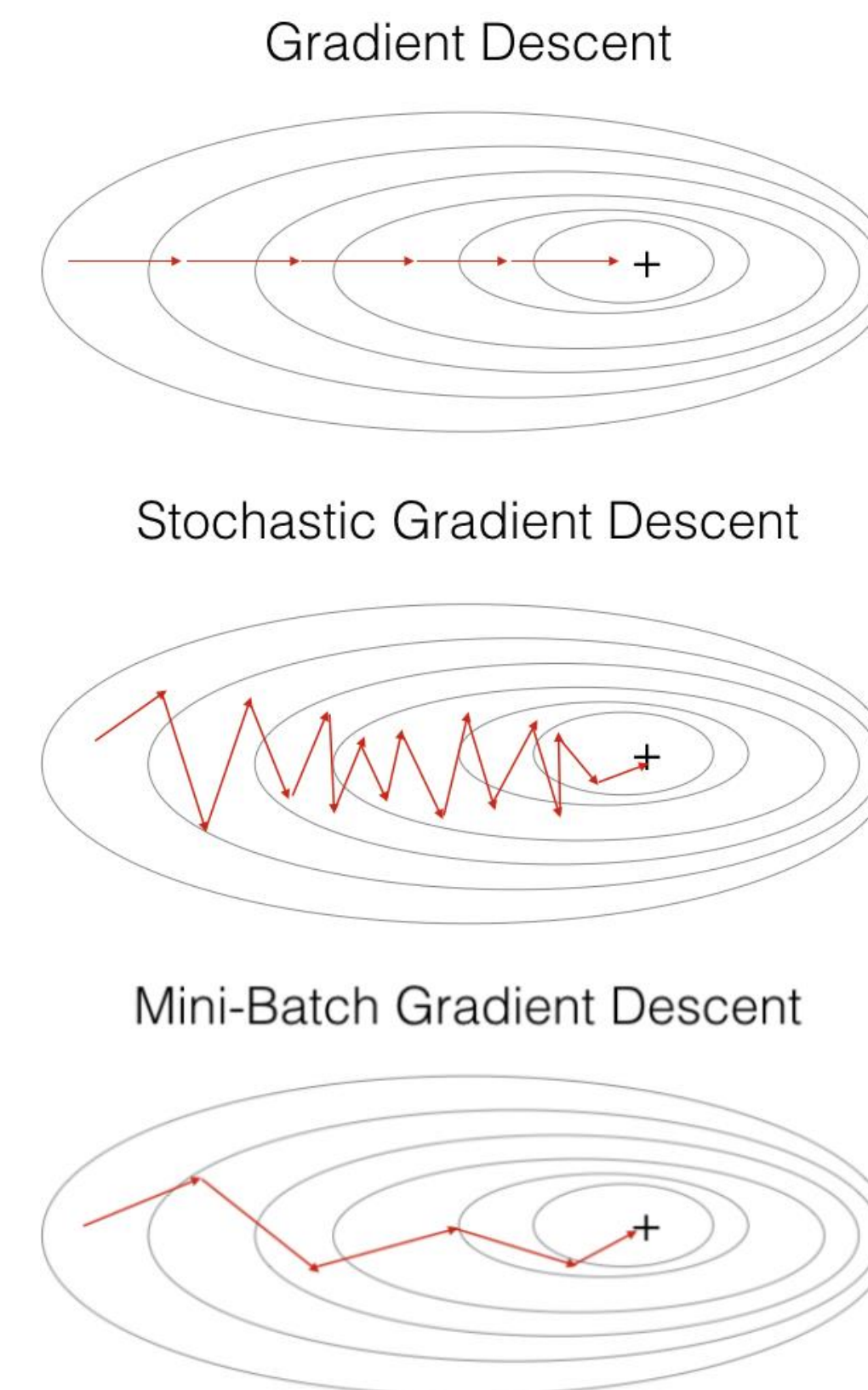
Η εποχή: 1 πέρασμα
πάνω από όλα τα
δεδομένα εκμάθησης





Gradient Descent

- Χρήση batch GD εάν το σύνολο εκμάθησης μπορεί να χωρέσει στη μνήμη
 - Ομαλή καμπύλη εκμάθησης
 - Αργή (1 ενημέρωση απαιτεί ολόκληρο το σύνολο εκπαίδευσης)
 - Μπορεί εύκολα να κολλήσει σε τοπικά ελάχιστα
- Χρησιμοποιήστε στοχαστικό GD για μεγάλα σύνολα κατάρτισης:
 - Θορυβώδης καμπύλη προπόνησης
 - Γρήγορη (1 ενημέρωση απαιτεί μία μόνο περίπτωση)
 - Εξακολουθεί να είναι επιρρεπής στο να κολλήσει σε τοπικά ελάχιστα, αλλά λιγότερο από το batch GD
- Στην πράξη:
 - Χρησιμοποιήστε mini-batch GD ρυθμίζοντας το μέγεθος παρτίδας ανάλογα με το πρόβλημα (π.χ. μέγεθος εμφάνισης) και τις προδιαγραφές του υπολογιστή έτσι ώστε να μπορεί να χωρέσει στη μνήμη
 - Μέγεθος παρτίδας: $2^0, 2^1, 2^2, \dots, m$



[ΠΗΓΗ](#)





Αποτελεσματικότητα της οπισθοδιάδοσης

Οπισθοδιάδοση:

- Χρειάζεται $O(W)$ λειτουργίας όπου W = αριθμός παραμέτρων

Αριθμητική διαφοροποίηση:

- Μέθοδος πεπερασμένων διαφορών:

$$\frac{\partial L(\theta_0, \dots, \theta_i, \dots, \theta_k)}{\partial \theta_i} \approx \frac{L(\theta_0, \dots, \theta_i + \varepsilon, \dots, \theta_k) - L(\theta_0, \dots, \theta_i - \varepsilon, \dots, \theta_k)}{2\varepsilon} \quad \varepsilon = 10^{-4}$$

- Ανάγκη να διαταράξει κάθε παράμετρο και να υπολογίσει την απώλεια: $O(W^2)$ πράξεις
- Μπορεί να χρησιμοποιηθεί για να ελέγξει αν η υλοποίηση της οπισθοδιάδοσης είναι σωστή





Διανυσματοποιημένη υλοποίηση

Διάδοση προς τα εμπρός:

$$\mathbf{a}^{(1)} = \mathbf{x}$$

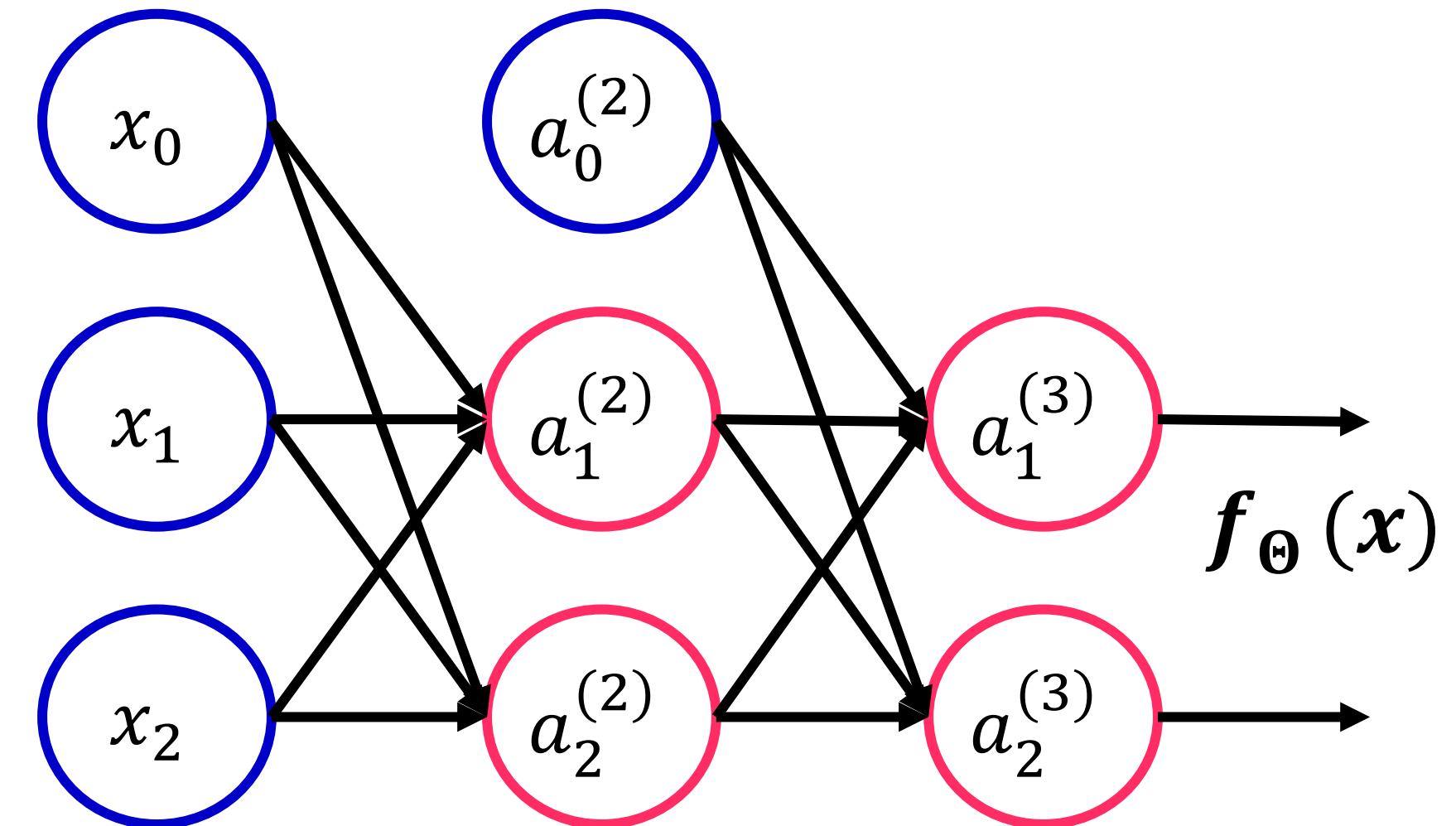
$$\mathbf{z}^{(2)} = \Theta^{(1)} \mathbf{a}^{(1)}$$

$$\mathbf{a}^{(2)} = \varphi(\mathbf{z}^{(2)})$$

Πρόσθεσε $a_0^{(2)} = 1$

$$\mathbf{z}^{(3)} = \Theta^{(2)} \mathbf{a}^{(2)}$$

$$\mathbf{a}^{(3)} = \varphi(\mathbf{z}^{(3)}) = \mathbf{f}_{\Theta}(\mathbf{x})$$





Διανυσματοποιημένη υλοποίηση

Οπισθοδιάδοση:

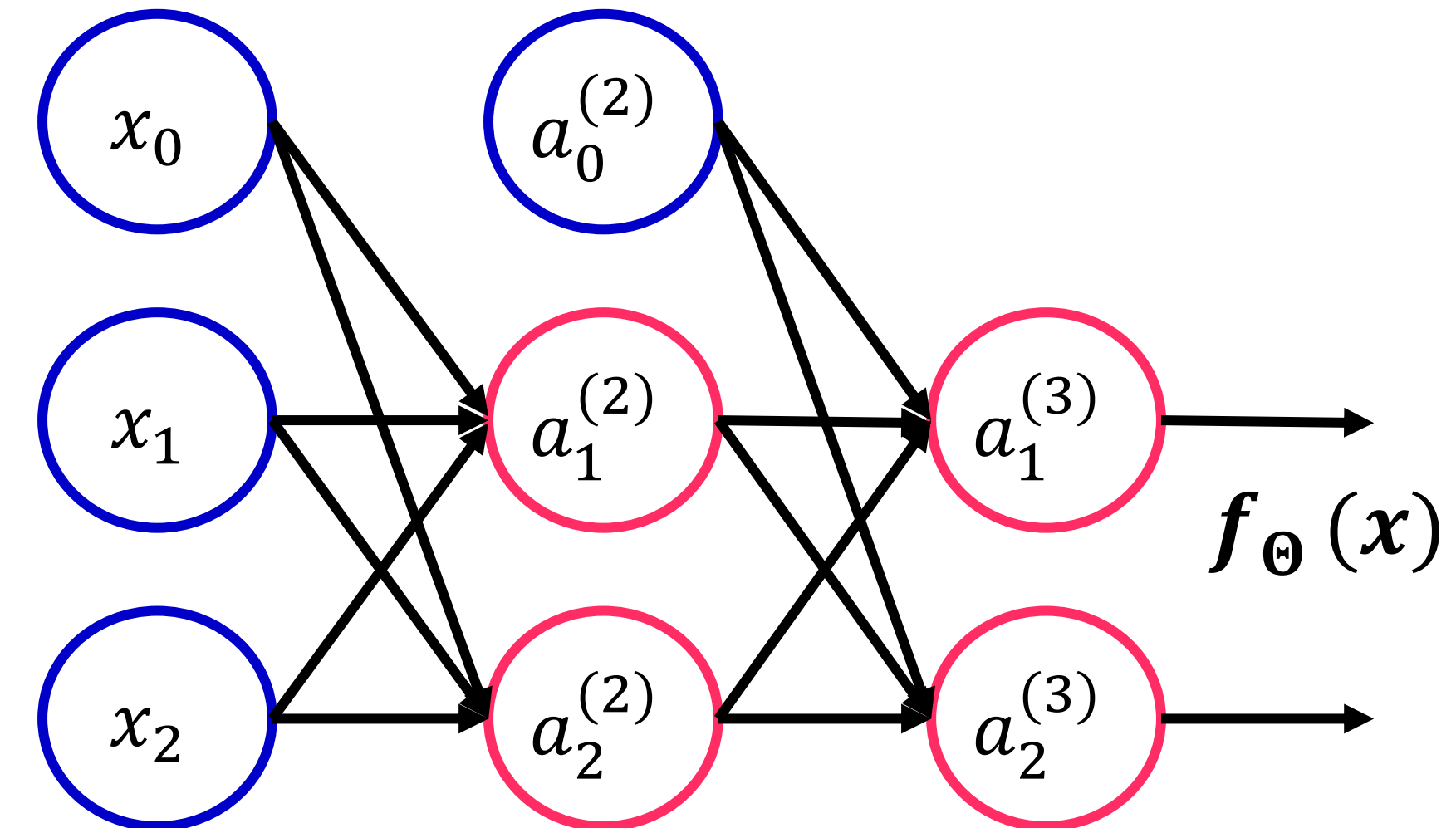
$$\delta^{(3)} = \mathbf{a}^{(3)} - \mathbf{y}$$

$$\delta^{(2)} = (\Theta^{(2)})^T \delta^{(3)} .* \boldsymbol{\varphi}'(\mathbf{z}^{(2)})$$

.* : element-wise multiplication

$$\boldsymbol{\varphi}'(\mathbf{z}^{(2)}) = \mathbf{a}^{(2)} .* (1 - \mathbf{a}^{(2)})$$

No $\delta^{(1)}$



$$\frac{\partial}{\partial \Theta_{j,i}^{(l)}} L(\Theta) = a_j^{(l)} \delta_i^{(l+1)}$$

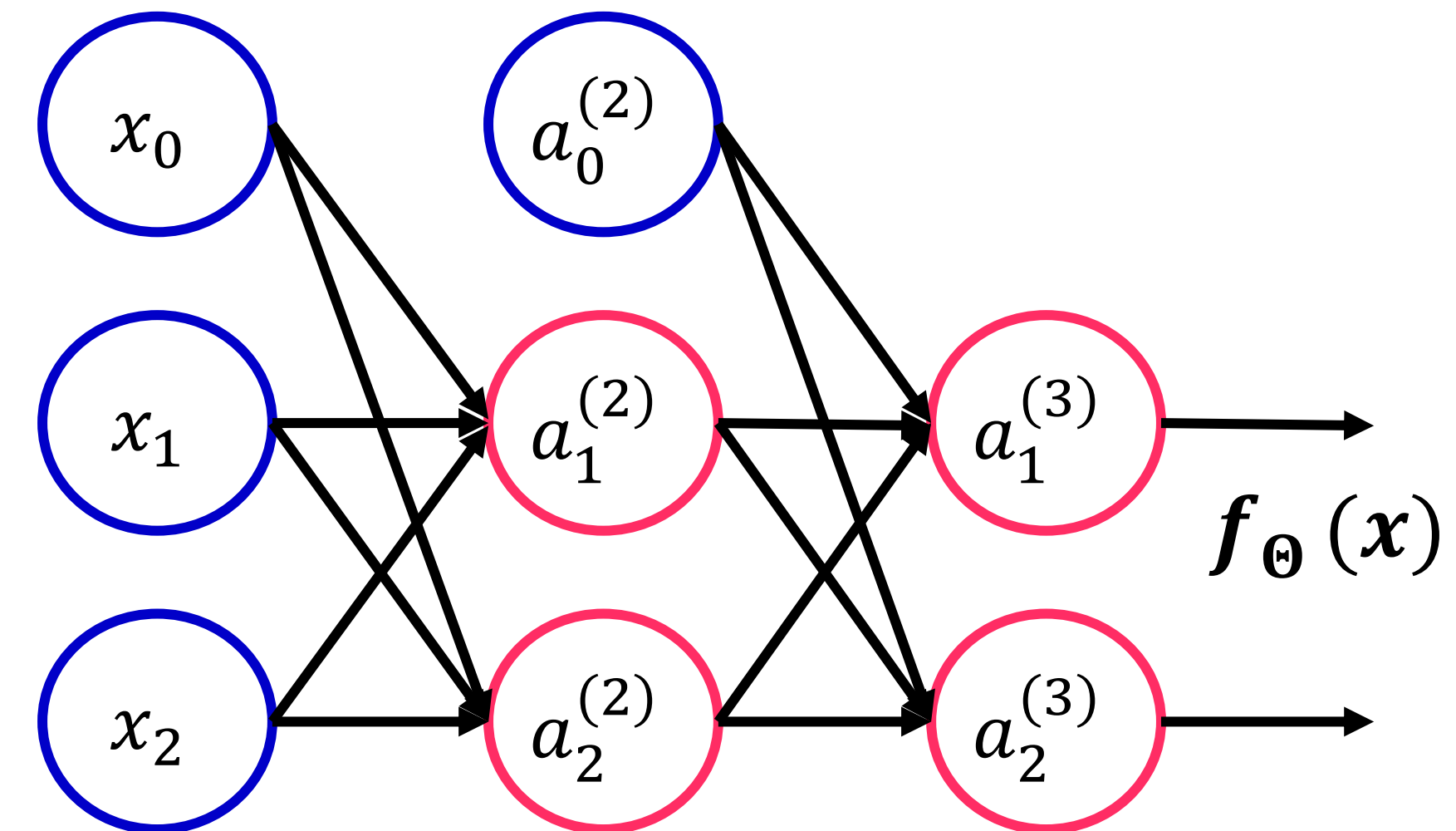




Τυχαία αρχικοποίηση

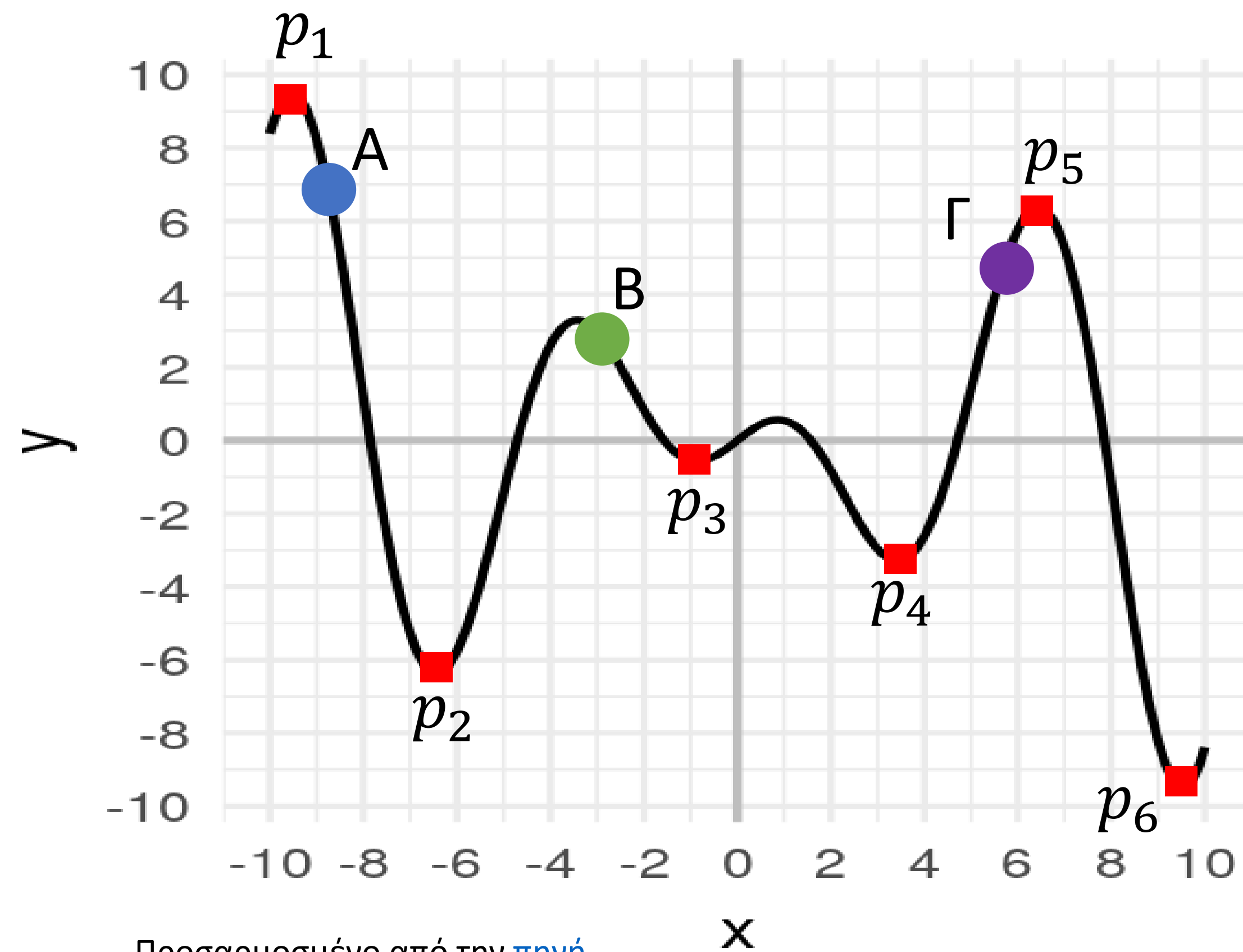
Πώς να αρχικοποιήσετε $\Theta_{j,i}^{(l)}$;

- Μηδενική αρχικοποίηση: $\Theta_{j,i}^{(l)} = 0$ για όλα τα i, j, l
 - Μετά από κάθε ενημέρωση, τα βάρη που έχουν την ίδια πηγή νευρώνων είναι πανομοιότυπα
 - Όλες οι κρυφές μονάδες υπολογίζουν την ίδια λειτουργία
 - Το δίκτυο δεν μαθαίνει καλά
- Τυχαία αρχικοποίηση: $\Theta_{j,i}^{(l)} \sim U(-\epsilon, \epsilon)$
 - Π.χ. $\epsilon = 1$
 - Κάθε κόμβος υπολογίζει μια διαφορετική συνάρτηση
 - Το δίκτυο μαθαίνει πιο ενδιαφέρουσες λειτουργίες





Κουίζ



Πού θα καταλήξει η βαθμιδωτή κάθοδος αν χρησιμοποιεί ένα μικρό ποσοστό μάθησης και ξεκινώντας από:

- A: p_2
- B: p_3
- Γ: p_4

Δεν θα φτάσει ποτέ στο p_6 .

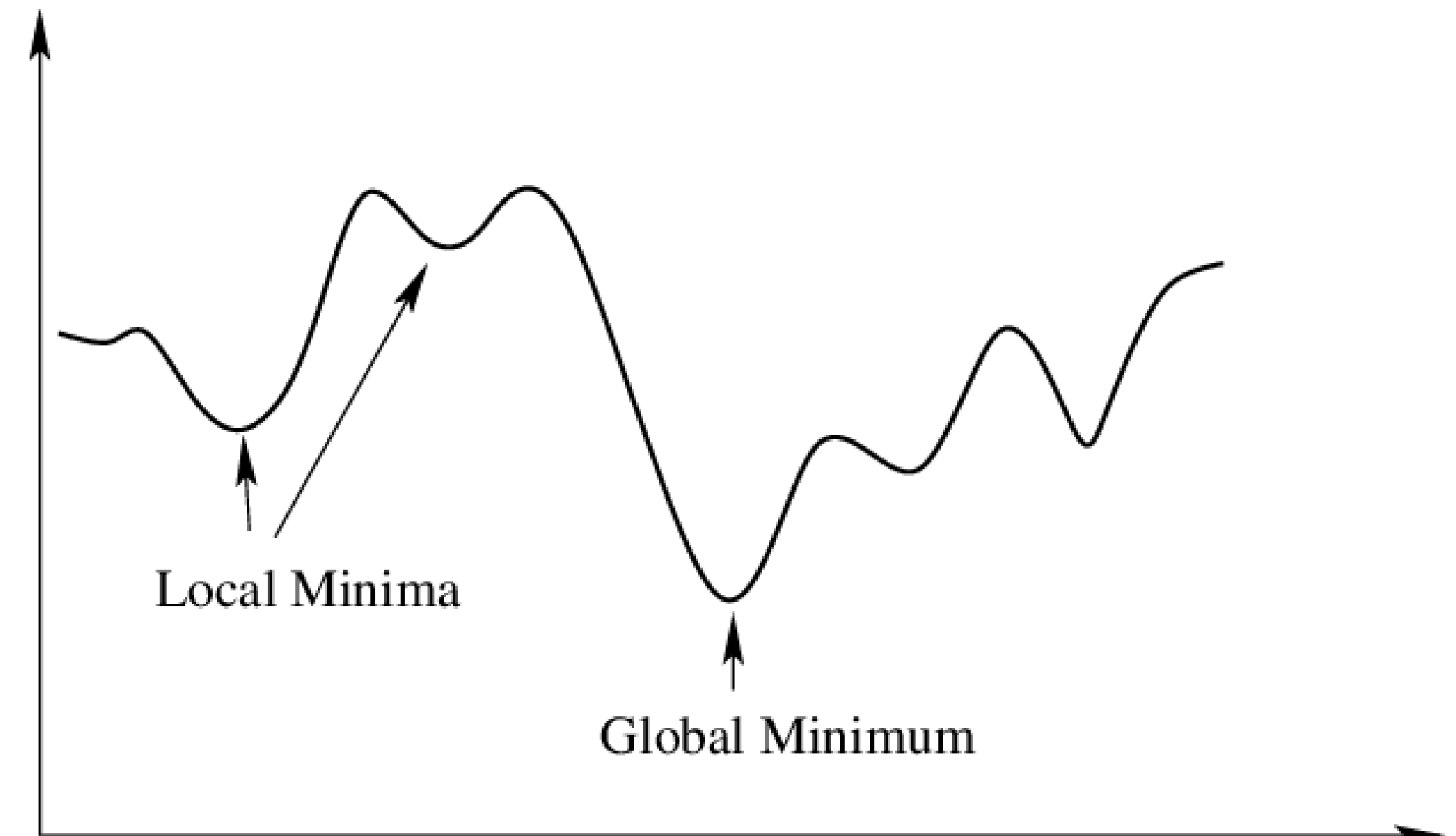




Προβλήματα με την SGD

- Να κολλήσει σε τοπικά ελάχιστα
- Αργή σύγκλιση ή απόκλιση

Ας αναθεωρήσουμε την
προσέγγιση...



[ΠΗΓΗ](#)





Stochastic gradient descent με ορμή

Gradient descent

$$\theta(t + 1) = \theta(t) - \eta \nabla L(\theta(t))$$

ρυθμός μάθησης

Gradient descent με ορμή

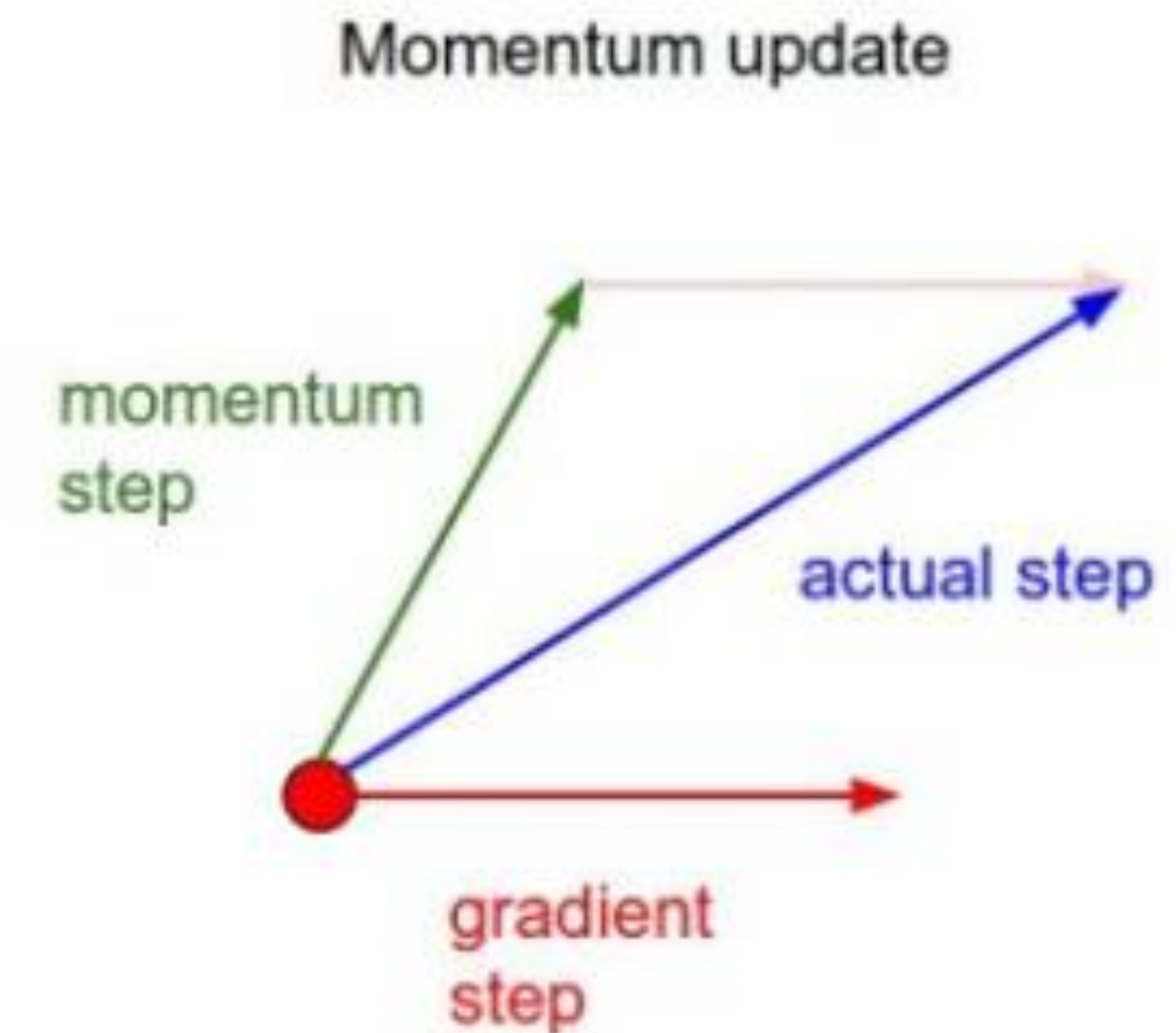
$$z(t + 1) = \mu z(t) + \nabla L(\theta(t))$$

$$\theta(t + 1) = \theta(t) - \eta z(t + 1)$$

συντελεστής ορμής

$$z(0) = 0$$

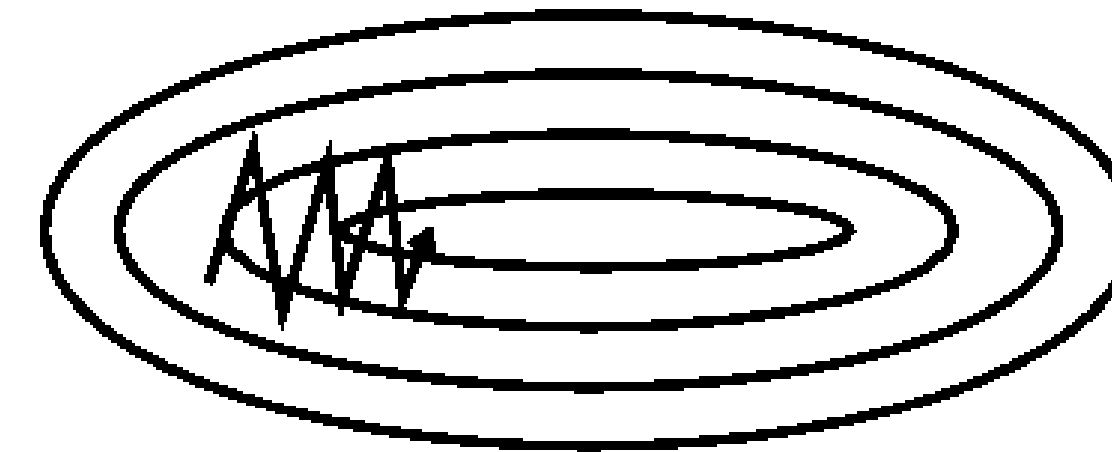
Αν $\mu = 0$ πάρουμε SGD



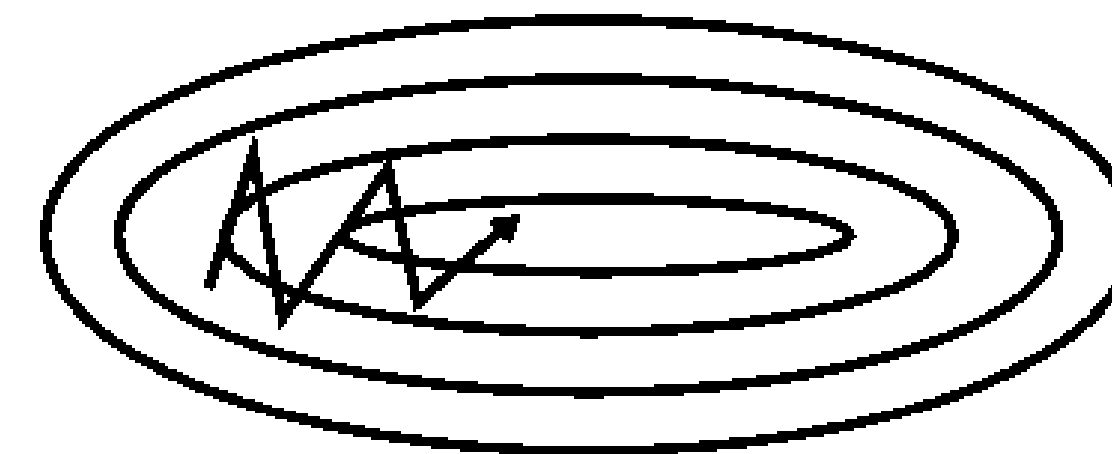


Stochastic gradient descent με ορμή

- Ορισμός ορμής: το είδος της μνήμης της προηγούμενης κατεύθυνσης που ακολουθούσαμε
- Μπορεί να επιταχύνει τη μάθηση
 - με την απόσβεση των ταλαντώσεων
 - με την αύξηση του αποτελεσματικού ποσοστού μάθησης σε περιοχές χαμηλής καμπυλότητας
- Μπορεί να ξεφύγει από το τοπικό βέλτιστο
- Χαρακτηριστική τιμή: $\mu = 0.9$



SGD

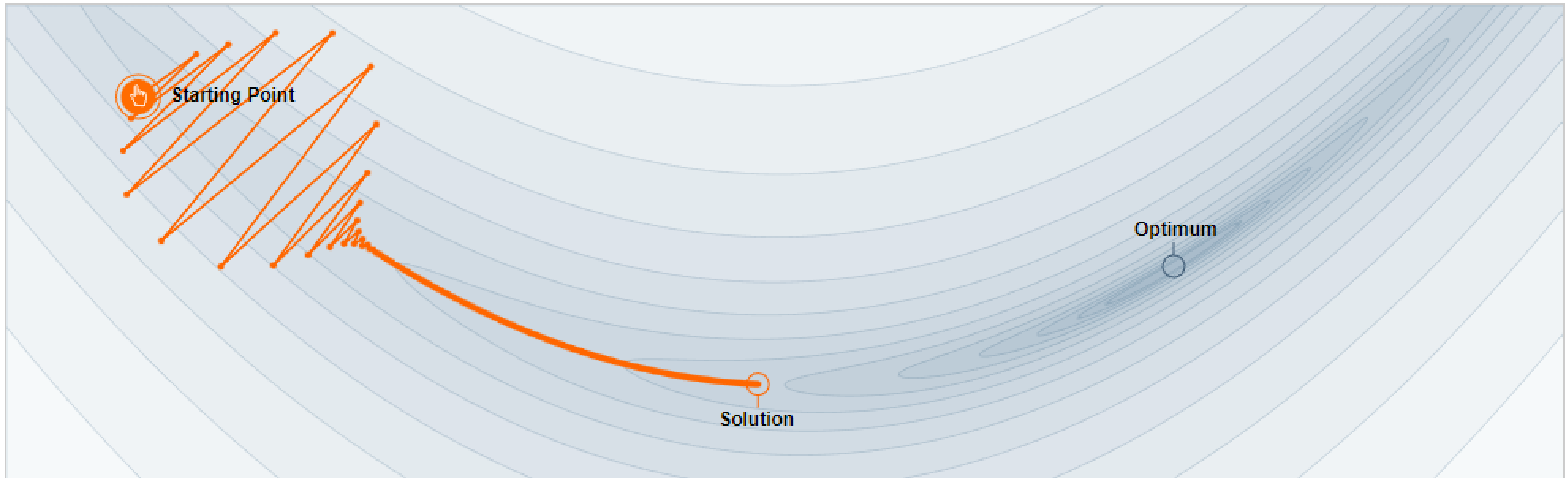


SGD με ορμή





SGD με ορμή



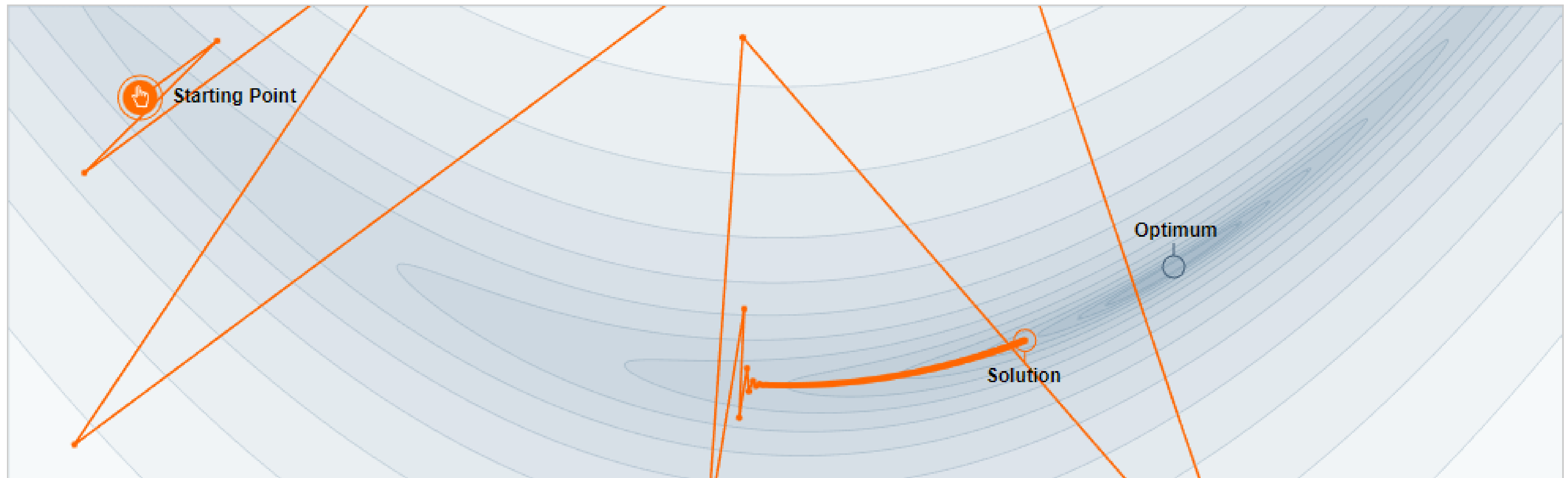
$$\eta = 0.003$$
$$\mu = 0.0$$

[ΠΗΓΗ](#)





SGD με ορμή



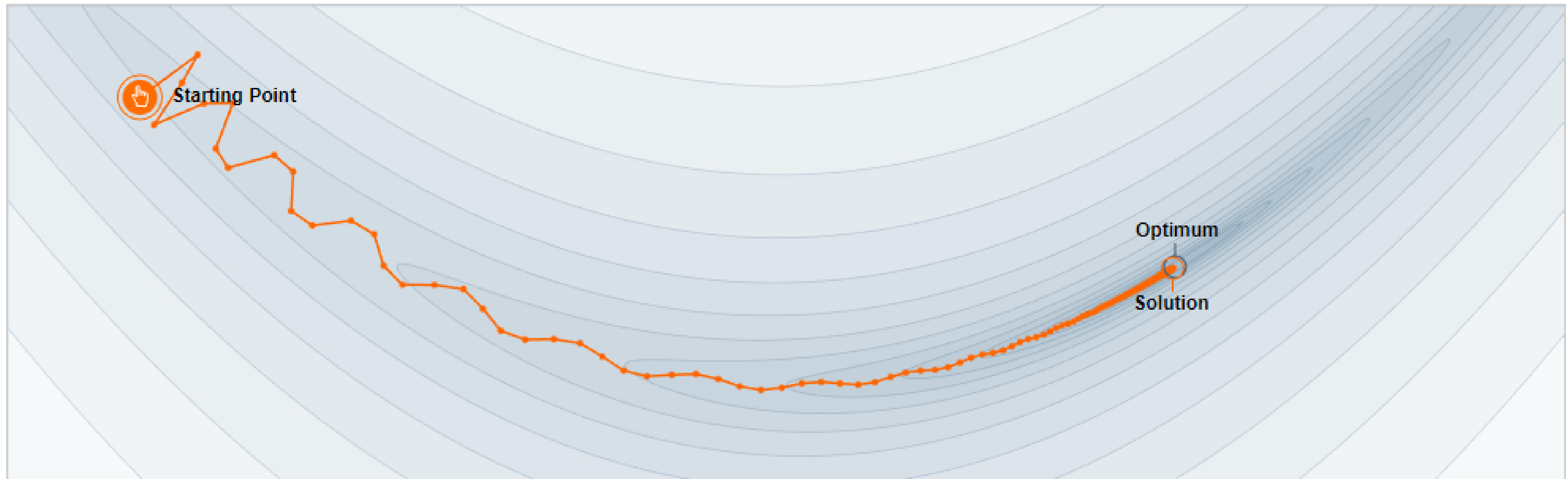
$$\eta = 0.004$$
$$\mu = 0.0$$

[ΠΗΓΗ](#)





SGD με ορμή



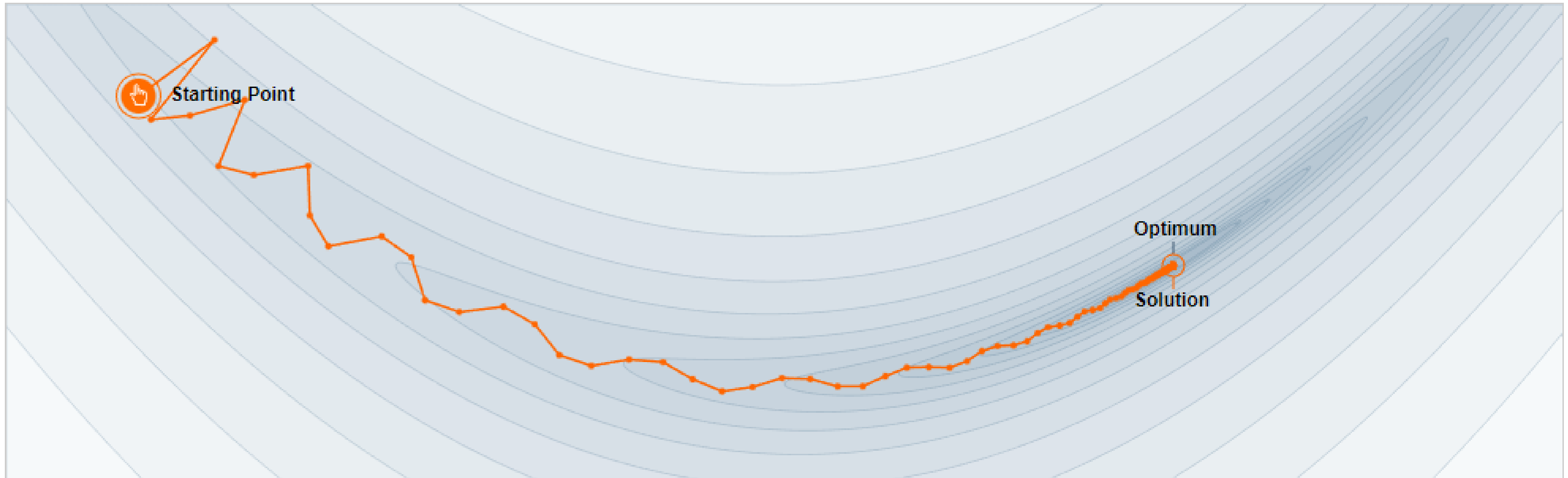
$$\eta = 0.003$$
$$\mu = 0.85$$

[ΠΗΓΗ](#)





SGD με ορμή



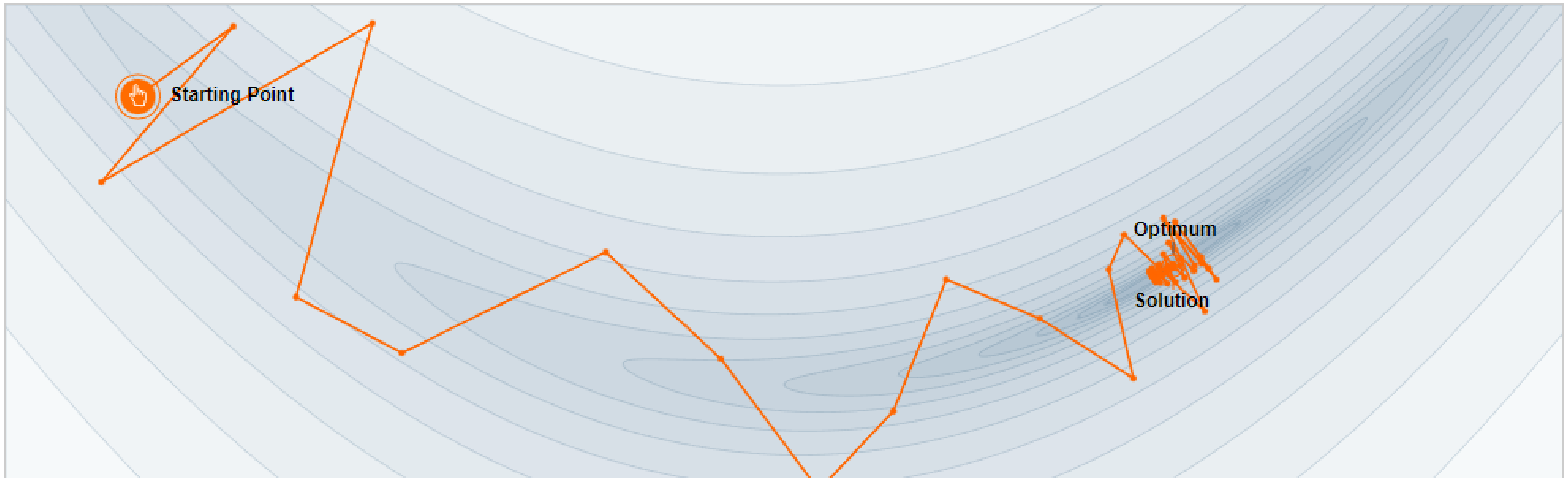
$$\eta = 0.004$$
$$\mu = 0.85$$

[ΠΗΓΗ](#)





SGD με ορμή



$$\eta = 0.005$$
$$\mu = 0.85$$

[ΠΗΓΗ](#)

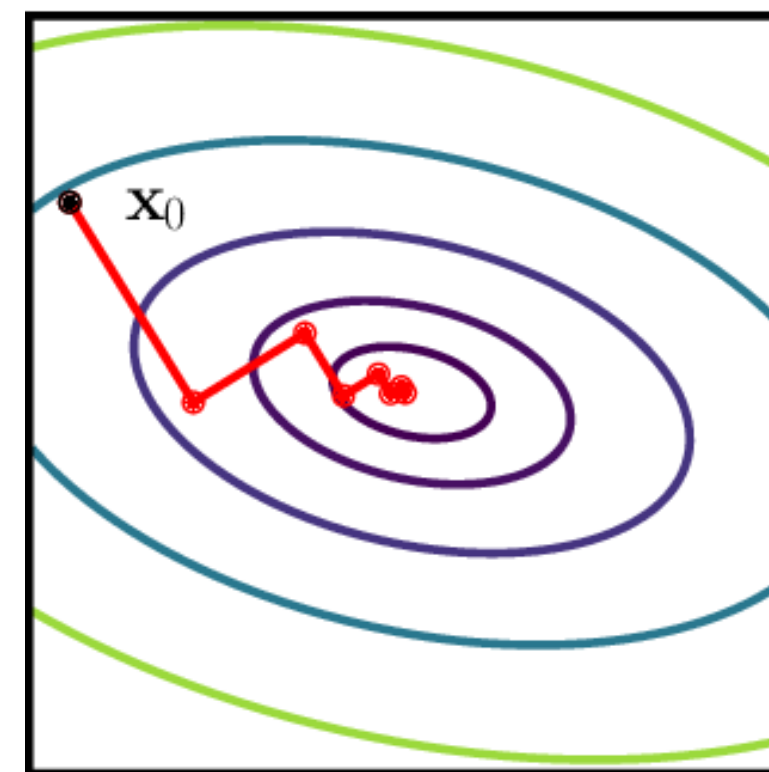




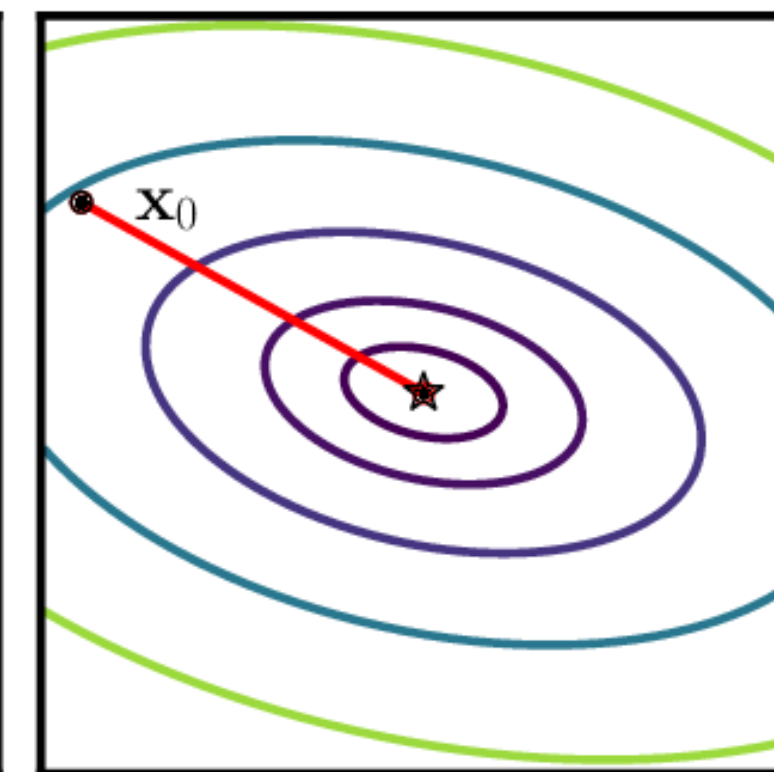
Προηγμένες μέθοδοι βελτιστοποίησης

- Η κάθοδος κλίσης είναι μια μέθοδος βελτιστοποίησης πρώτης τάξης
- Οι μέθοδοι δεύτερης τάξης χρησιμοποιούν πληροφορίες καμπυλότητας (2^η παράγωγο) να επιταχύνουν τη σύγκλιση στο βέλτιστο
 - Η μέθοδος του Νεύτωνα
 - Μέθοδος συζευγμένης κλίσης
 - Αλγόριθμος Levenberg-Marquardt
 - ...

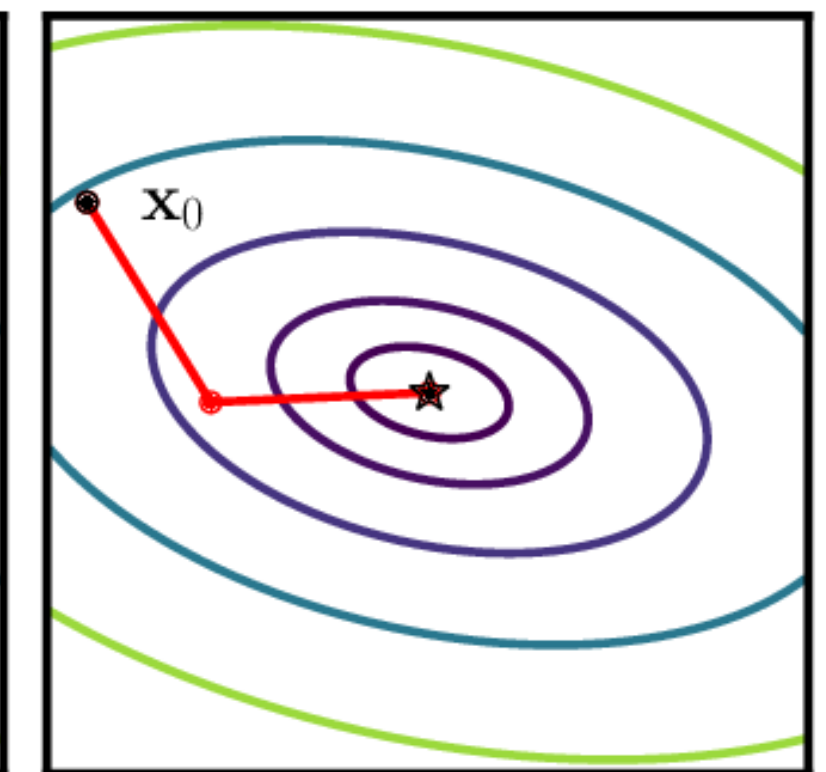
[ΠΗΓΗ](#)



(a) Steepest descent



(b) Newton's method



(c) Conjugate gradient





Προηγμένες μέθοδοι βελτιστοποίησης

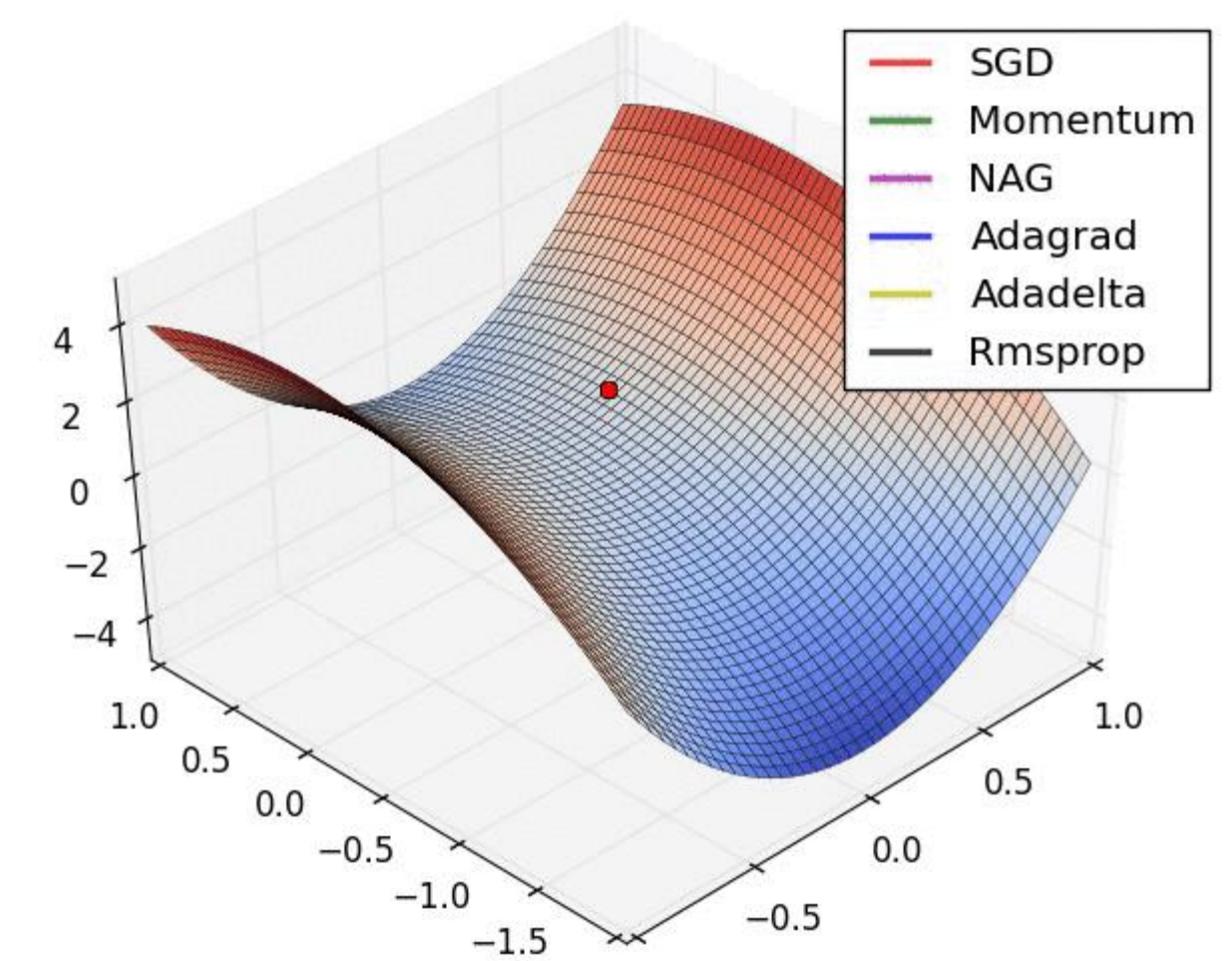
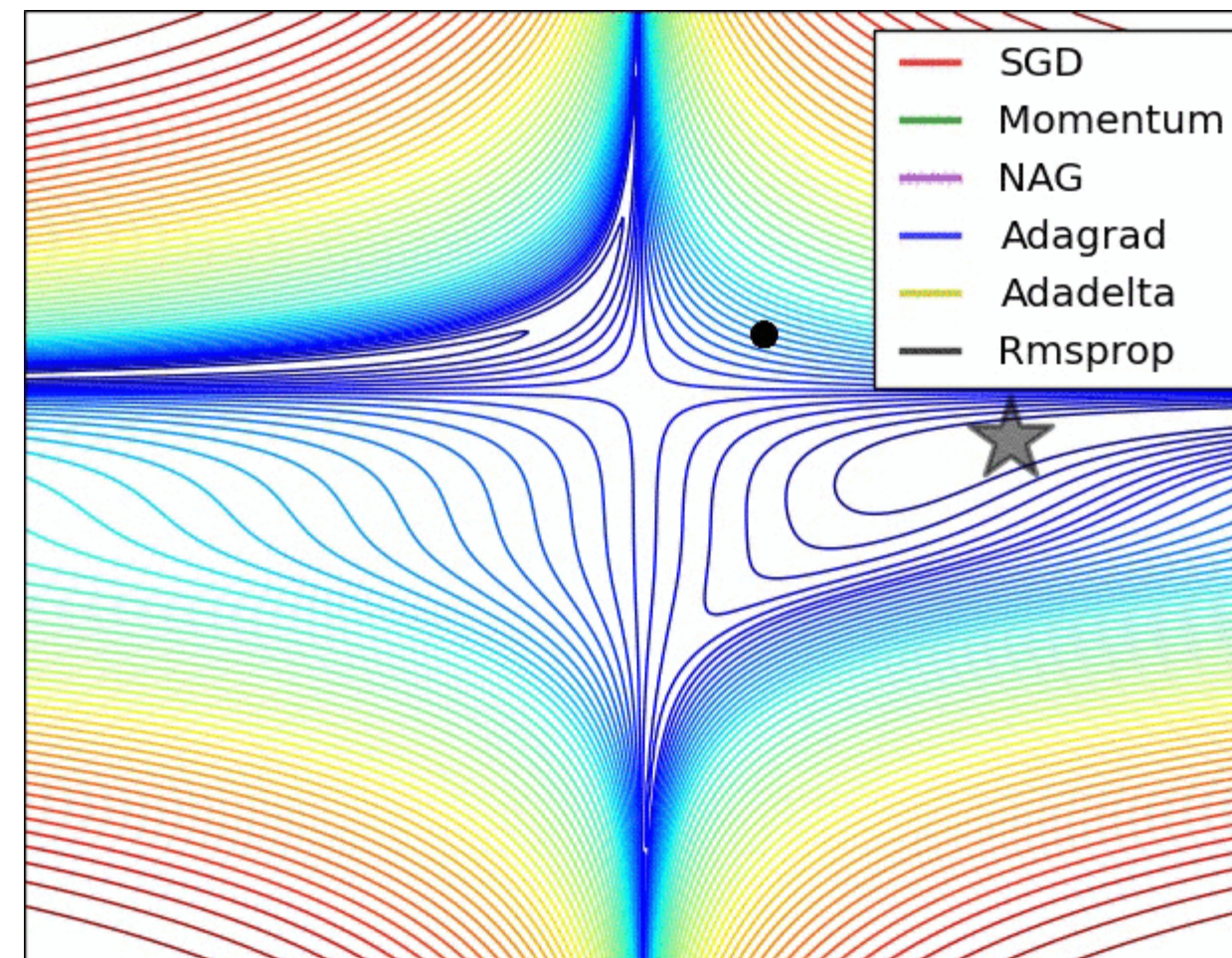
- Μέθοδοι δεύτερης τάξης
 - Πιο περίπλοκο και πιο δύσκολο να εφαρμοστεί και να συντονιστεί
 - Πιο ακριβά όσον αφορά το κόστος επανάληψης και τη μνήμη
 - Εξετάζουν την πρόσβαση στην πραγματική κλίση
 - Είναι πιο δύσκολο να τους κάνει να δουλέψουν με την ενημέρωση mini-batch
- Η σύγχρονη εκπαίδευση νευρωνικών δικτύων βασίζεται σε παραλλαγές της SGD





Προηγμένες μέθοδοι βελτιστοποίησης: Παραλλαγές SGD

- Απλό SQD
- SQD με ορμή
- Nesterov's accelerated gradient
- Adagrad
- Adadelta
- RMSProp
- Adam
- AdaMax
- Nadam
- ...

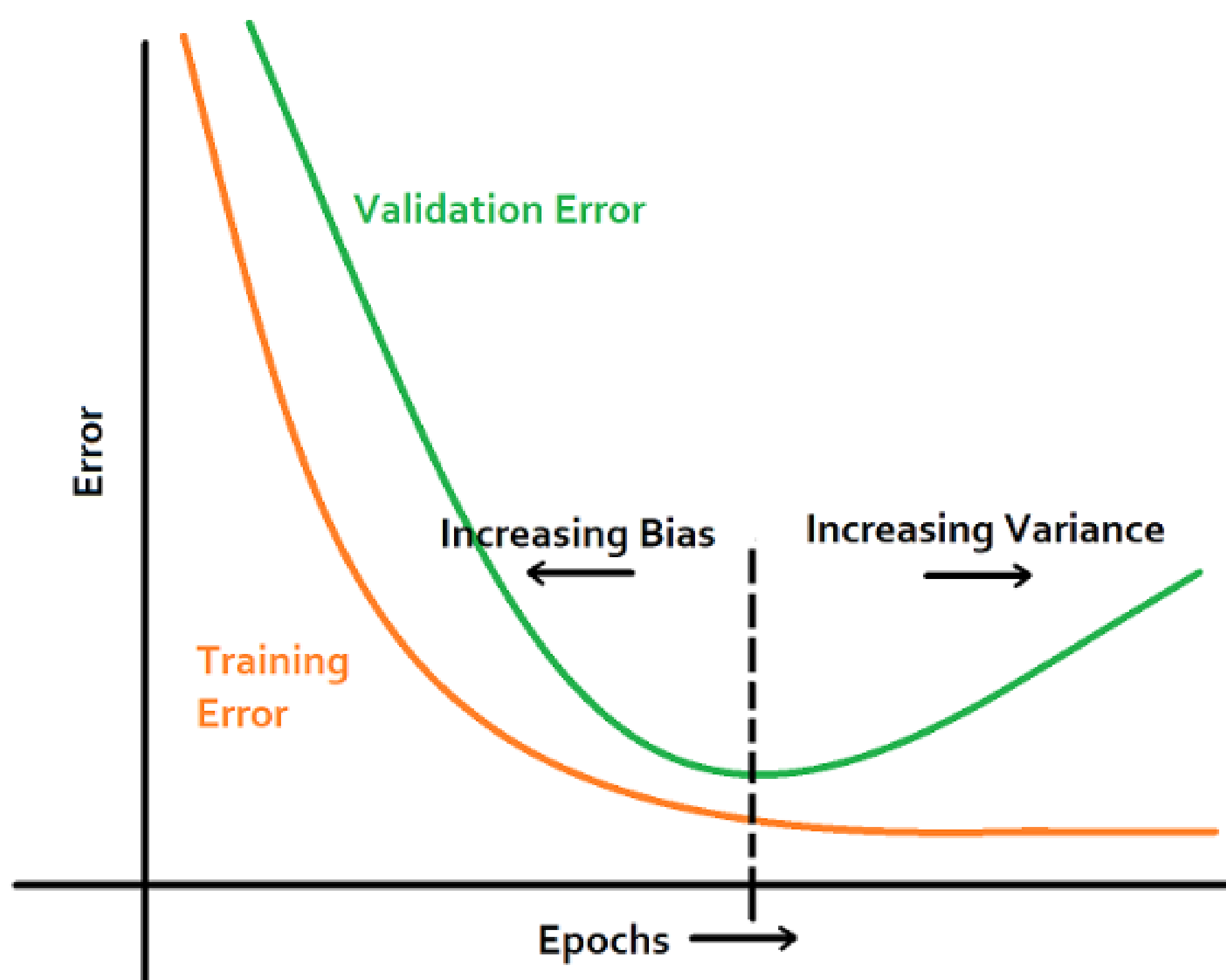


[ΠΗΓΗ](#)





Πρόωρη διακοπή



[ΠΗΓΗ](#)

- Τα ΝΔ είναι μοντέλα με μεγάλη ικανότητα μάθησης: επιρρεπής στην υπερπροσαρμογή

Πρόωρη διακοπή ως τακτοποίηση:

- Κατά την εκμάθηση, αξιολογήστε την απόδοση στο σύνολο επικύρωσης
- Αποθηκεύστε τα βάρη δικτύου εάν το σφάλμα επικύρωσης μειώνεται
- Όταν το σφάλμα επικύρωσης αρχίσει να αυξάνεται (συνεπώς) να σταματήσει και να επιστρέψει το δίκτυο με τα αποθηκευμένα βάρη





Ρύθμιση υπερπαραμέτρων

Υπερπαραμέτροι μοντέλου:

- Αριθμός κρυφών στρωμάτων
- Αριθμός μονάδων ανά στρώμα
- Λειτουργίες ενεργοποίησης ανά στρώμα (ανά νευρώνα;)
- Παράμετροι λειτουργίας ενεργοποίησης

Υπερπαραμέτροι αλγορίθμων εκμάθησης:

- Ποσοστό μάθησης
- Συντελεστής ορμής
- Φθορά ποσοστού μάθησης
- Παράμετροι αρχικοποίησης
- Κανονικοποίηση
- ...

Η επιλογή μοντέλου χρησιμοποιείται πριν από την αξιολόγηση των συνδυασμών υπερπαραμέτρων χρησιμοποιώντας το σύνολο επικύρωσης.





Ensembles

- Τα ΝΔ έχουν πολλές παραμέτρους: το σφάλμα έχει πολλές διαστάσεις με πολλά βέλτιστα
- Το Gradient descent είναι μια θορυβώδης διαδικασία που δεν είναι εγγυημένη ότι συγκλίνει στο ίδιο τοπικό βέλτιστο για κάθε εκτέλεση.
- Μπορούμε να εκπαιδεύσουμε πολλαπλά ανεξάρτητα δίκτυα και να μετρήσουμε τις επιδόσεις τους
- Μπορούμε να χρησιμοποιήσουμε bagging ή boosting ή πιο προηγμένες μεθόδους (π.χ. dropout) και να λάβουμε χαμηλότερα σφάλματα επικύρωσης





Αυτόματη διαφοροποίηση

- Για διαφορετικές συναρτήσεις κόστους ή συναρτήσεις ενεργοποίησης πρέπει να γράψουμε ρητό κώδικα που υπολογίζει την παράγωγο τους
- Για μια πιο ευέλικτη αρχιτεκτονική νευρωνικού δικτύου, πρέπει να γράψουμε σαφή κώδικα για τη διάδοση προς τα εμπρός και προς τα πίσω.

Αυτόματη διαφοροποίηση:

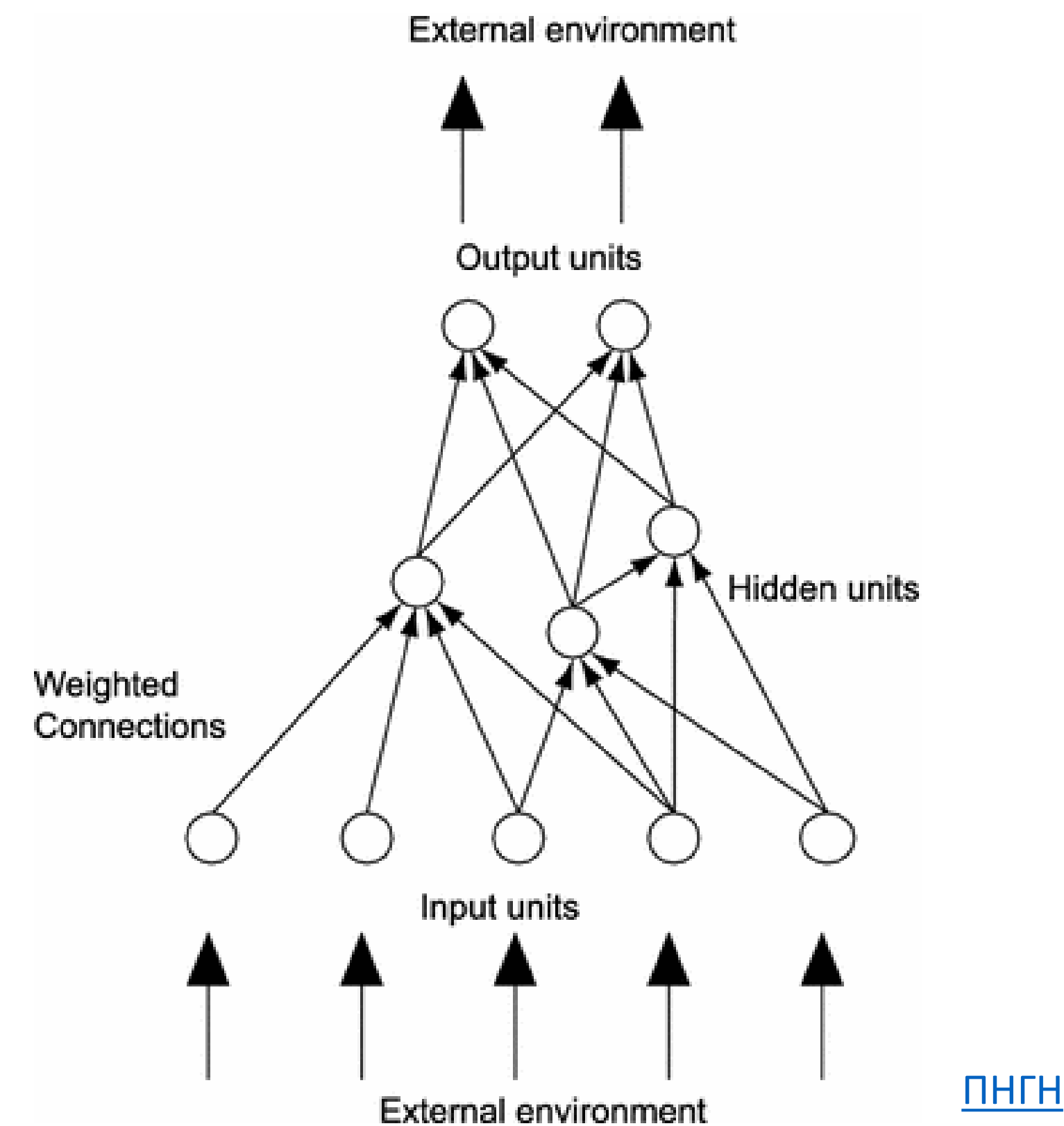
- Ένας γενικευμένος τρόπος χρήσης του κανόνα της αλυσίδας με διάφορες λειτουργίες και λειτουργίες
- Ουσιαστικά: Η οπισθοδιάδοση υλοποιείται για εμάς
- Βιβλιοθήκες: [TensorFlow](#), [PyTorch](#), [JAX](#)





Μαθαίνοντας τη δομή του δικτύου

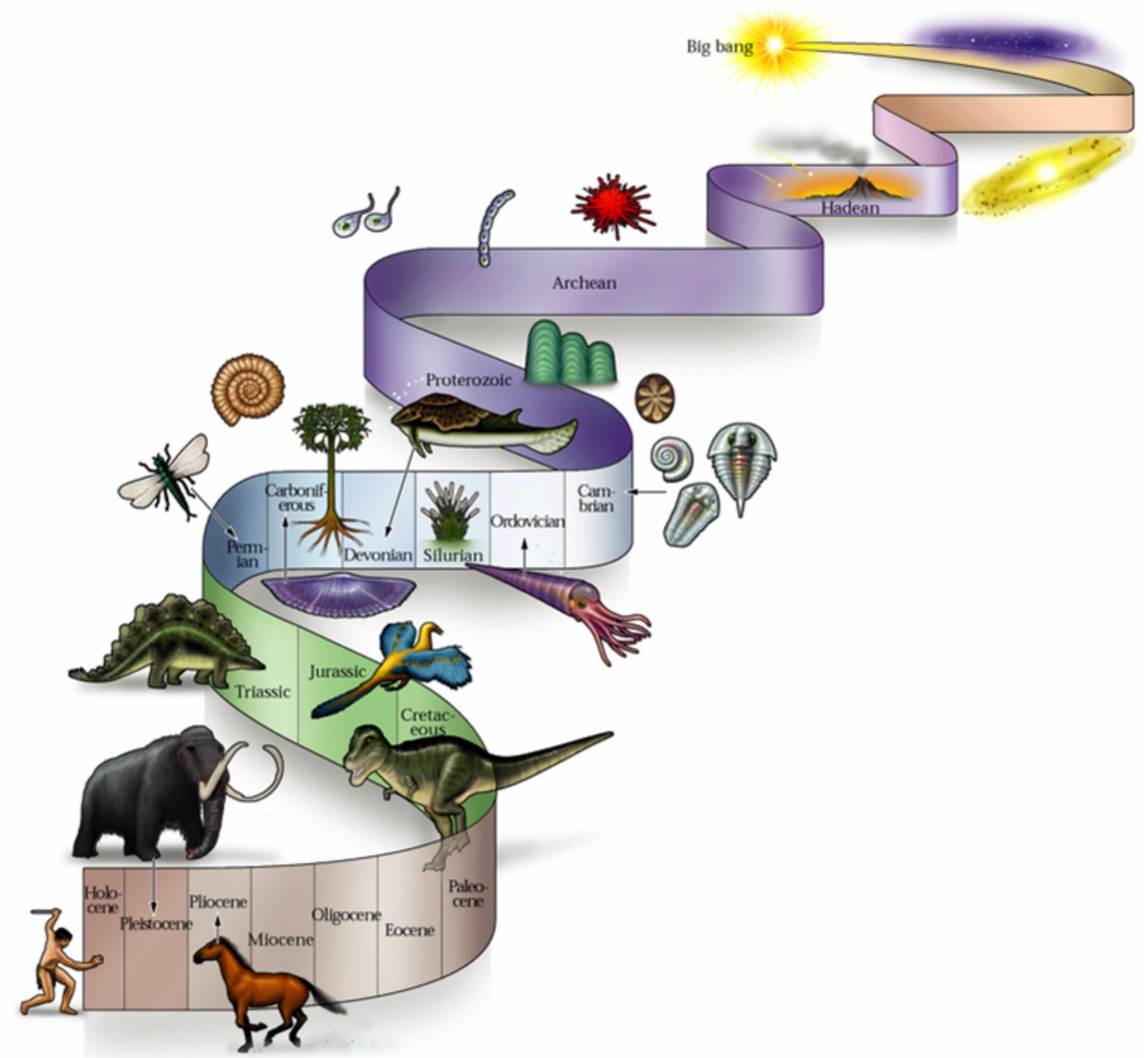
- Τα νευρωνικά δίκτυα είναι συνθέσεις λειτουργιών
- Η σύνθεση αυτών των συναρτήσεων με πολυεπίπεδο τρόπο τροφοδοτείται (καθολικό θεώρημα προσέγγισης), ωστόσο, υπάρχουν πιο αποδοτικοί τρόποι σύνθεσης τους για ένα δεδομένο πρόβλημα.
- Η εκμάθηση της συνδεσιμότητας δεν μπορεί να διατυπωθεί (εύκολα) ως στόχος βελτιστοποίησης βάσει gradient.
 - Χρειάζονται μεθόδους χωρίς gradient
 - Αναζήτηση πλέγματος ή τυχαία αναζήτηση μπορεί να είναι αναποτελεσματική
- **Neuroevolution**: χρήση **εξελικτικών αλγορίθμων** για το σχεδιασμό/βελτιστοποίηση των νευρωνικών δικτύων





Εξελικτικά συστήματα

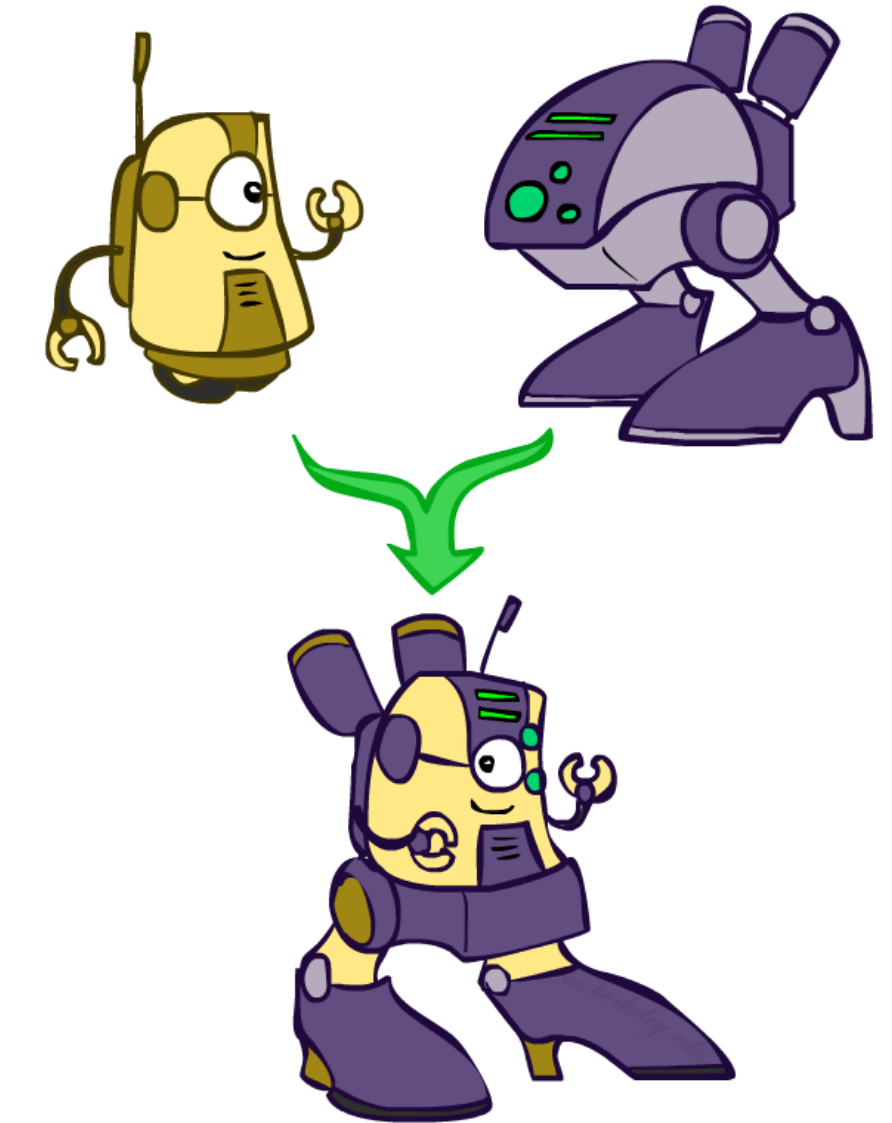
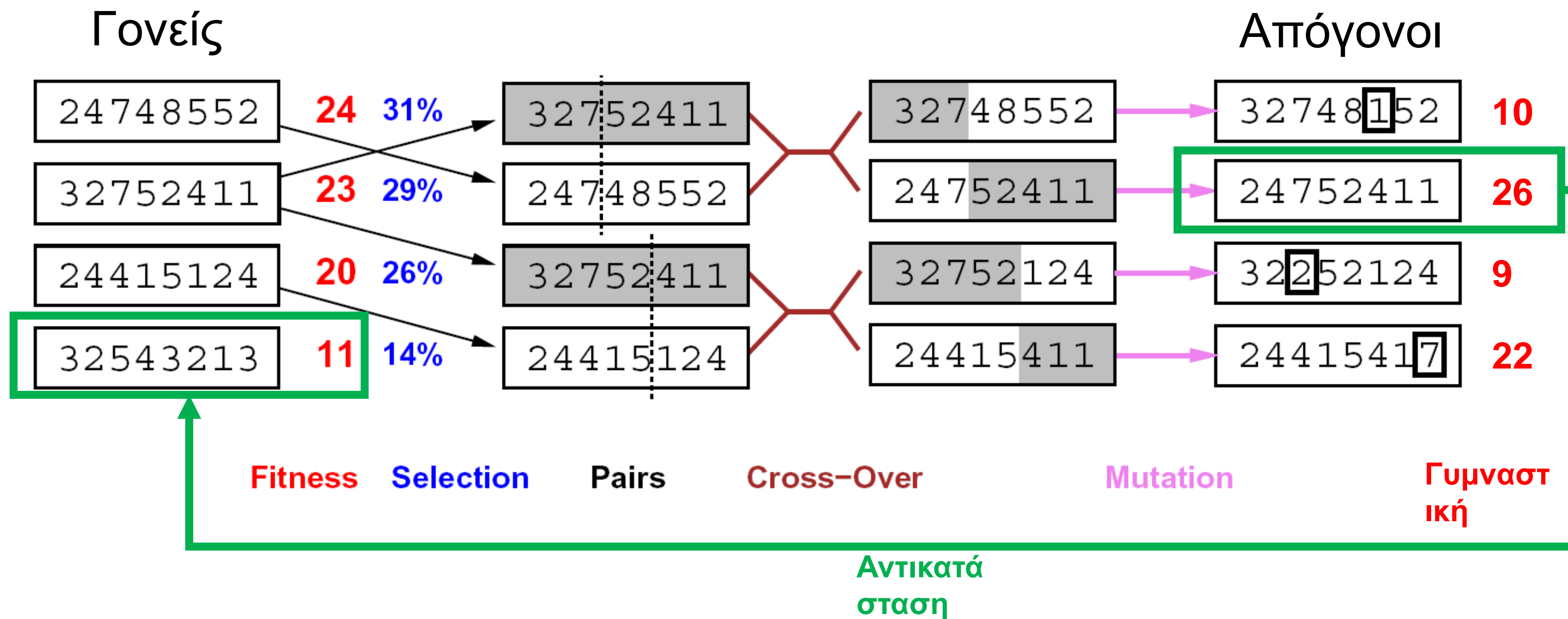
- Όλα τα βιολογικά συστήματα είναι το αποτέλεσμα μιας εξελικτικής διαδικασίας
- Τα βιολογικά συστήματα είναι:
 - στιβαρά
 - περίπλοκα
 - προσαρμοστικά
- Οι εξελικτικοί αλγόριθμοι είναι αφαιρέσεις φυσικών εξελικτικών διαδικασιών με στόχο την αυτόματη εξεύρεση λύσεων σε σύνθετα προβλήματα.
 - Π.χ. ως γενικοί βελτιστοποιητές





[ΠΗΓΗ](#)

Γενετικοί Αλγόριθμοι

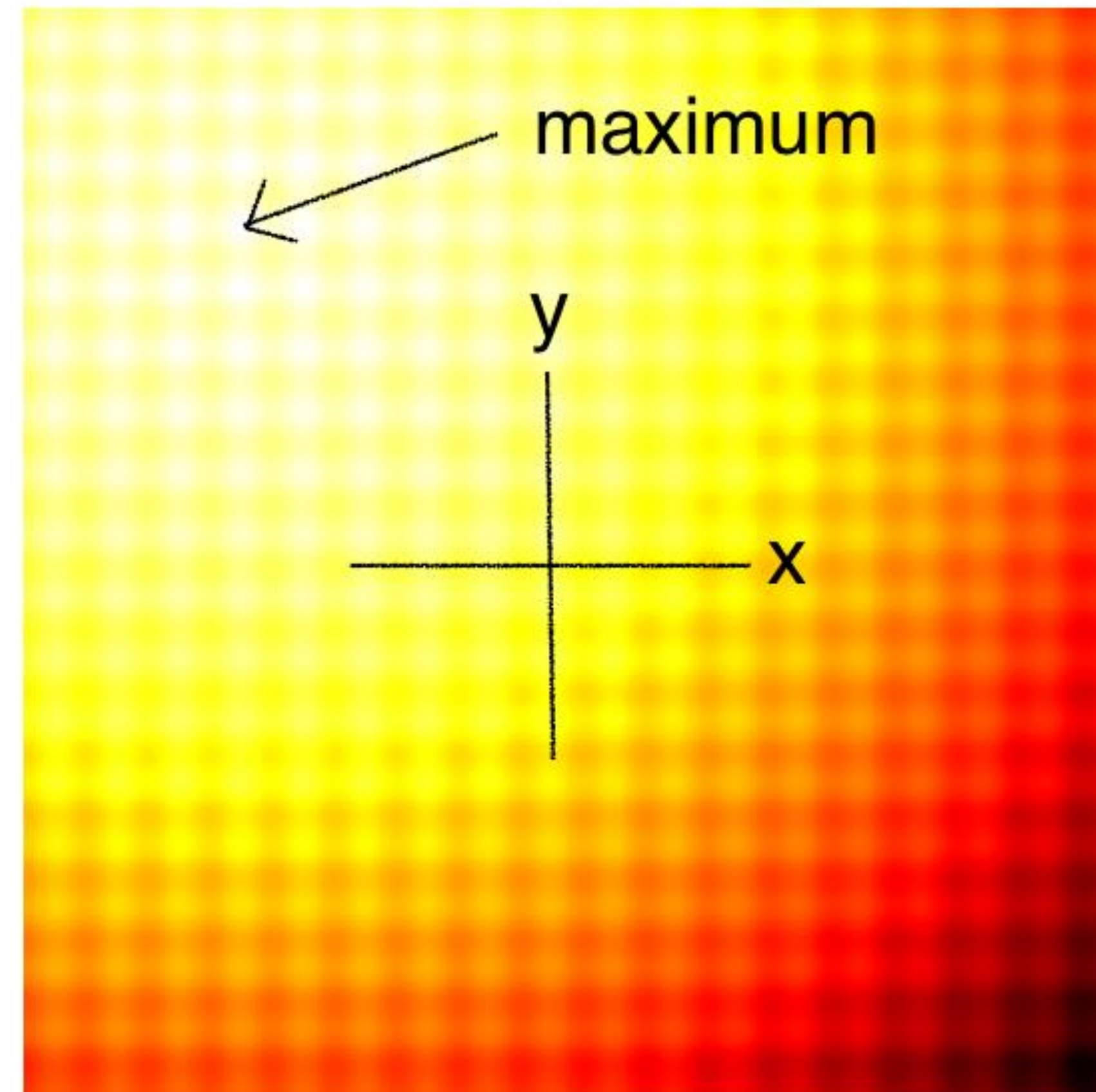
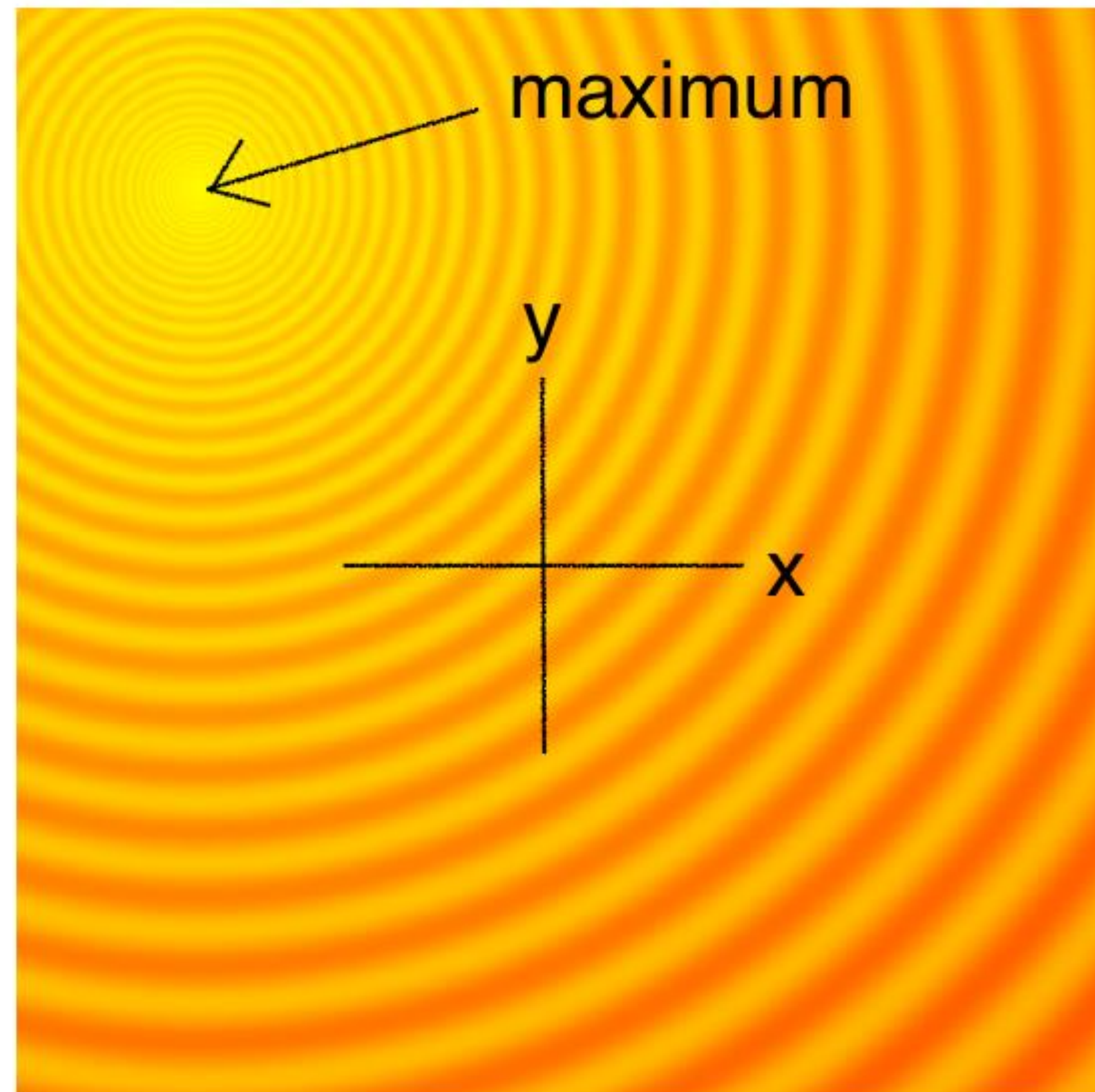


- Μεταφορά φυσικής επιλογής
 - Οι γονείς παράγουν απογόνους μέσω διασταύρωσης και μετάλλαξης γονιδίων
 - Εάν ο απόγονος έχει χαρακτηριστικά που είναι επωφελή για το περιβάλλον του (καλύτερη φυσική κατάσταση), επιβιώνει, έτσι, επιτρέπεται να διαδώσει τα γονίδιά του στην επόμενη γενιά.





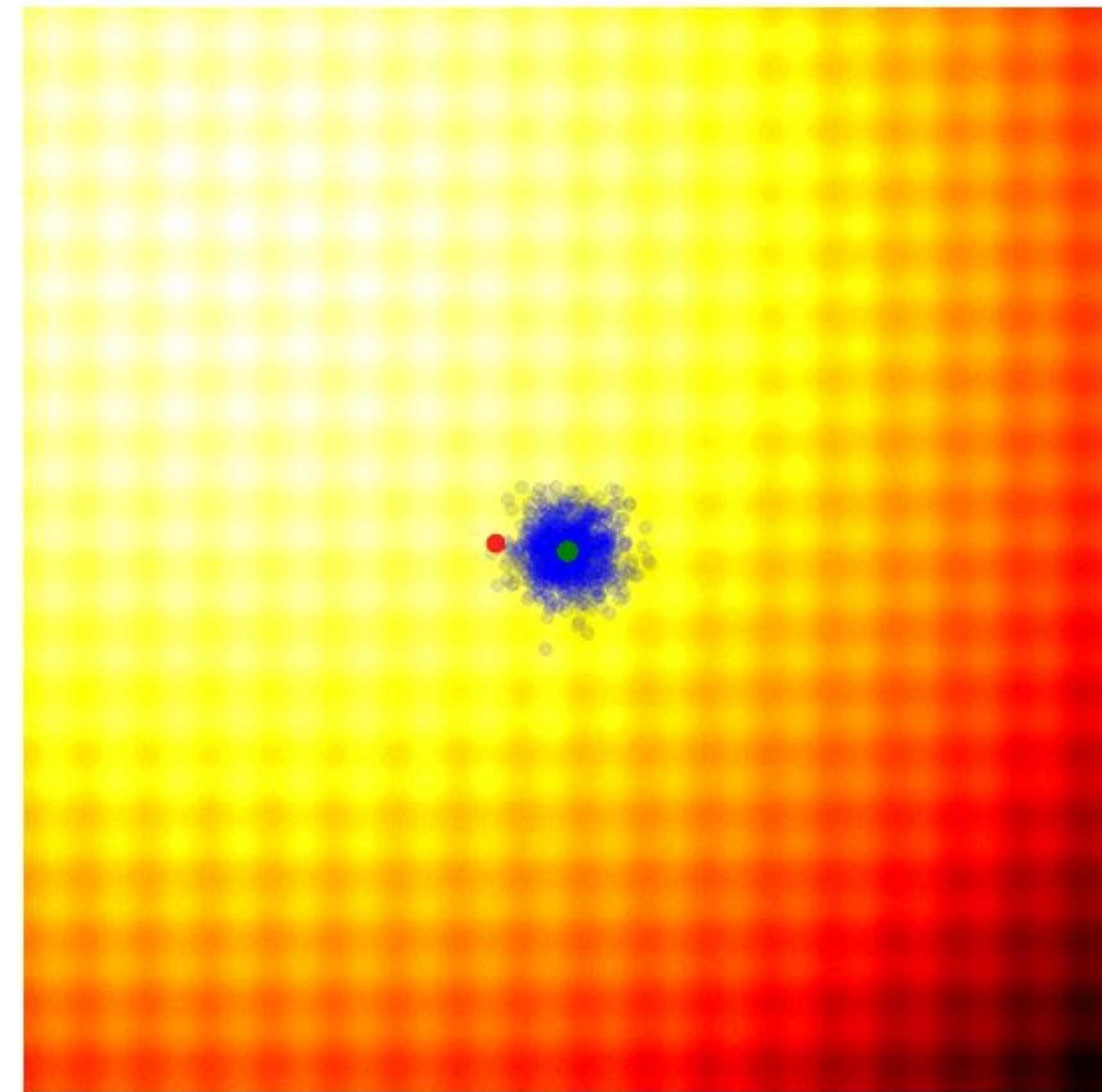
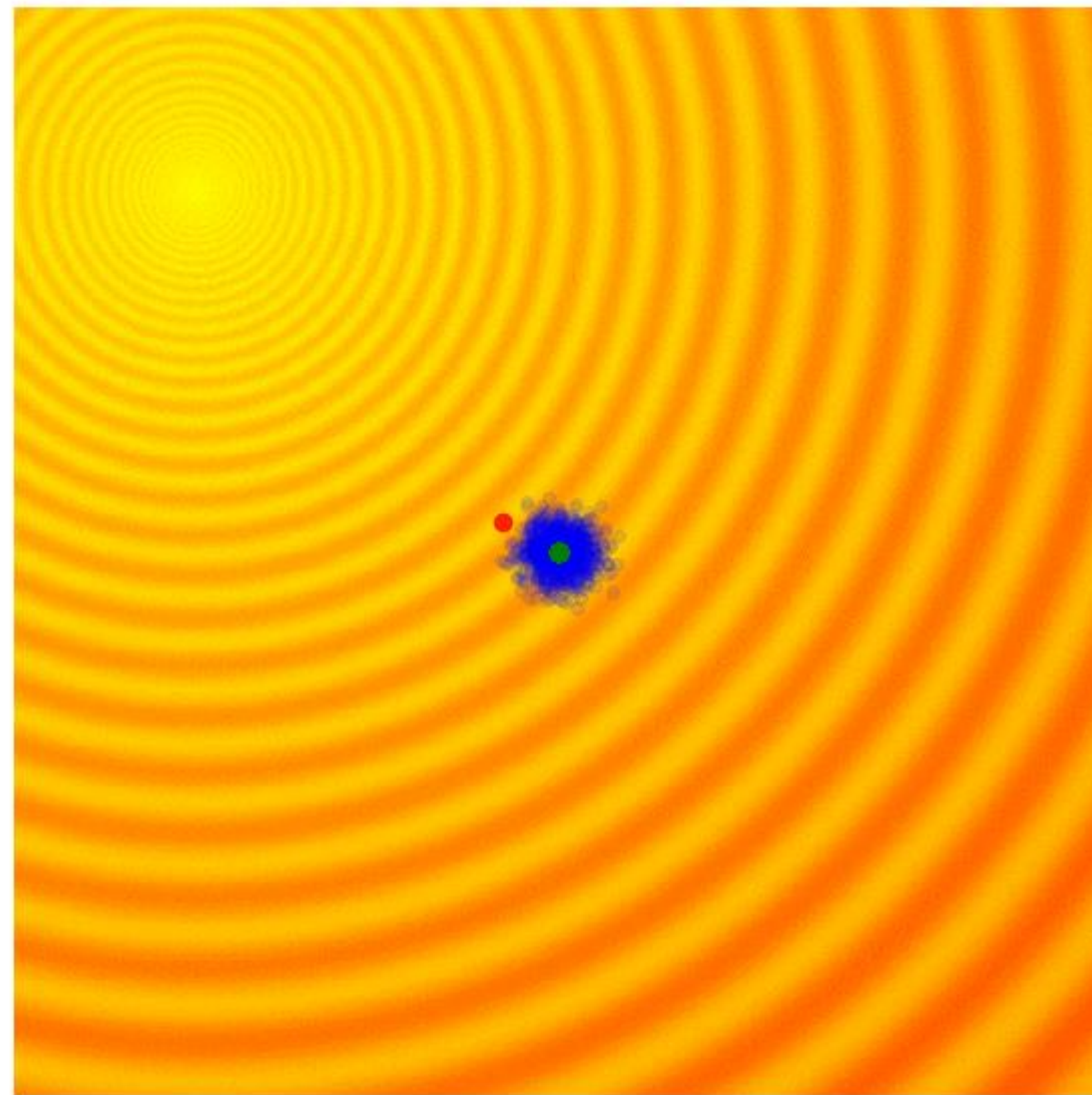
Εξελικτικοί Αλγόριθμοι σε Δράση



[ΠΗΓΗ](#)



Γενετικοί Αλγόριθμοι

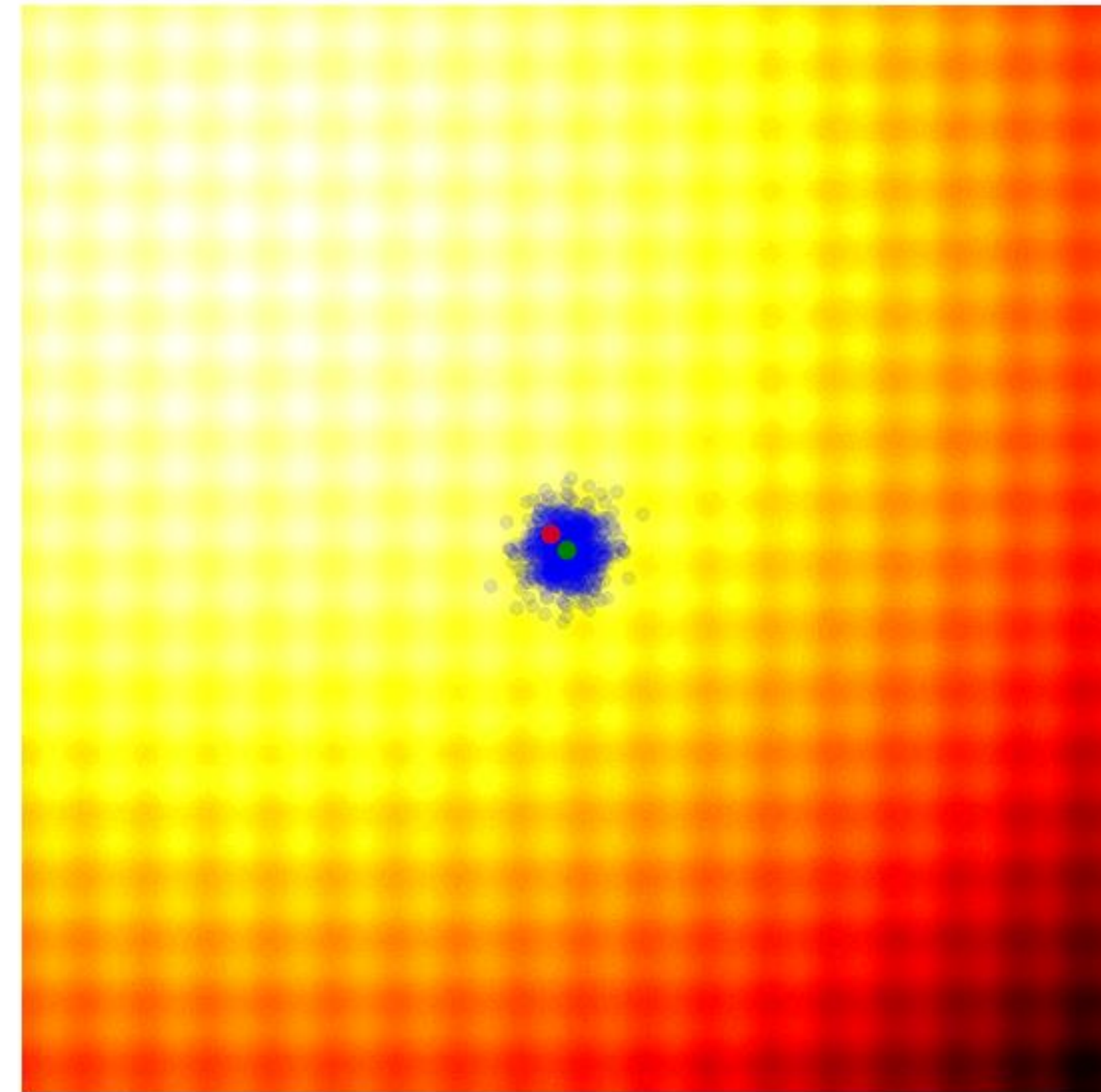
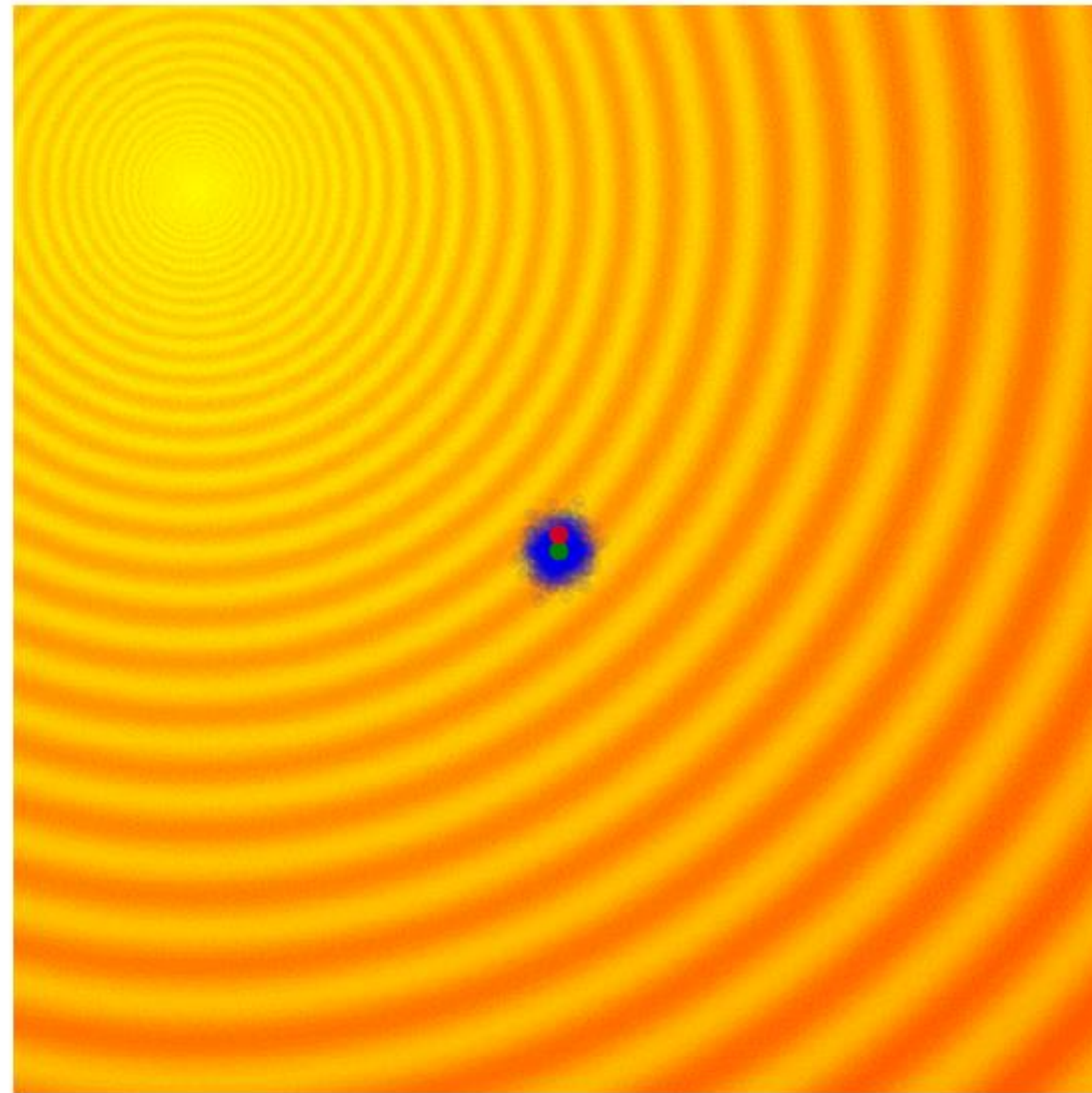


[ΠΗΓΗ](#)





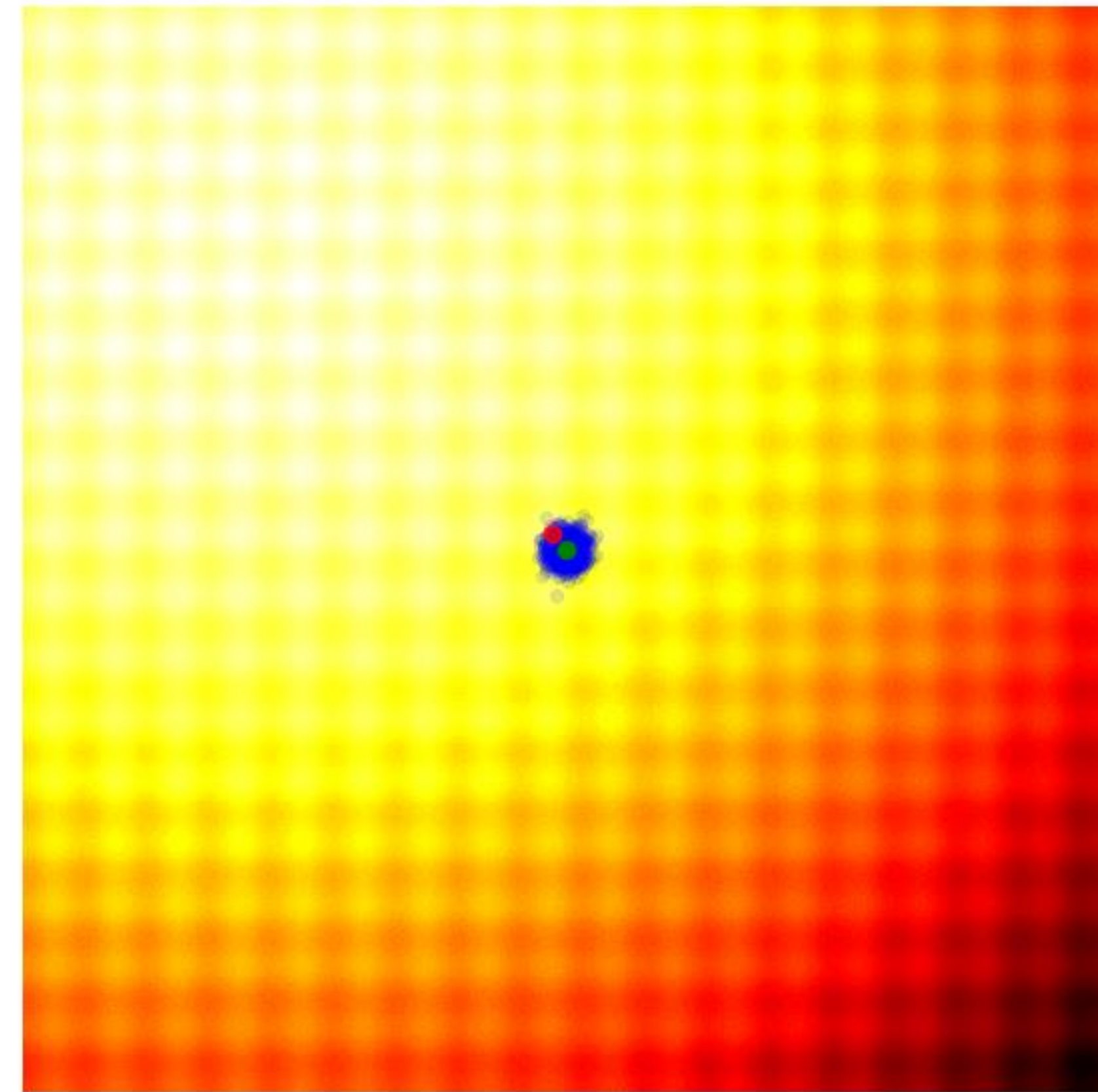
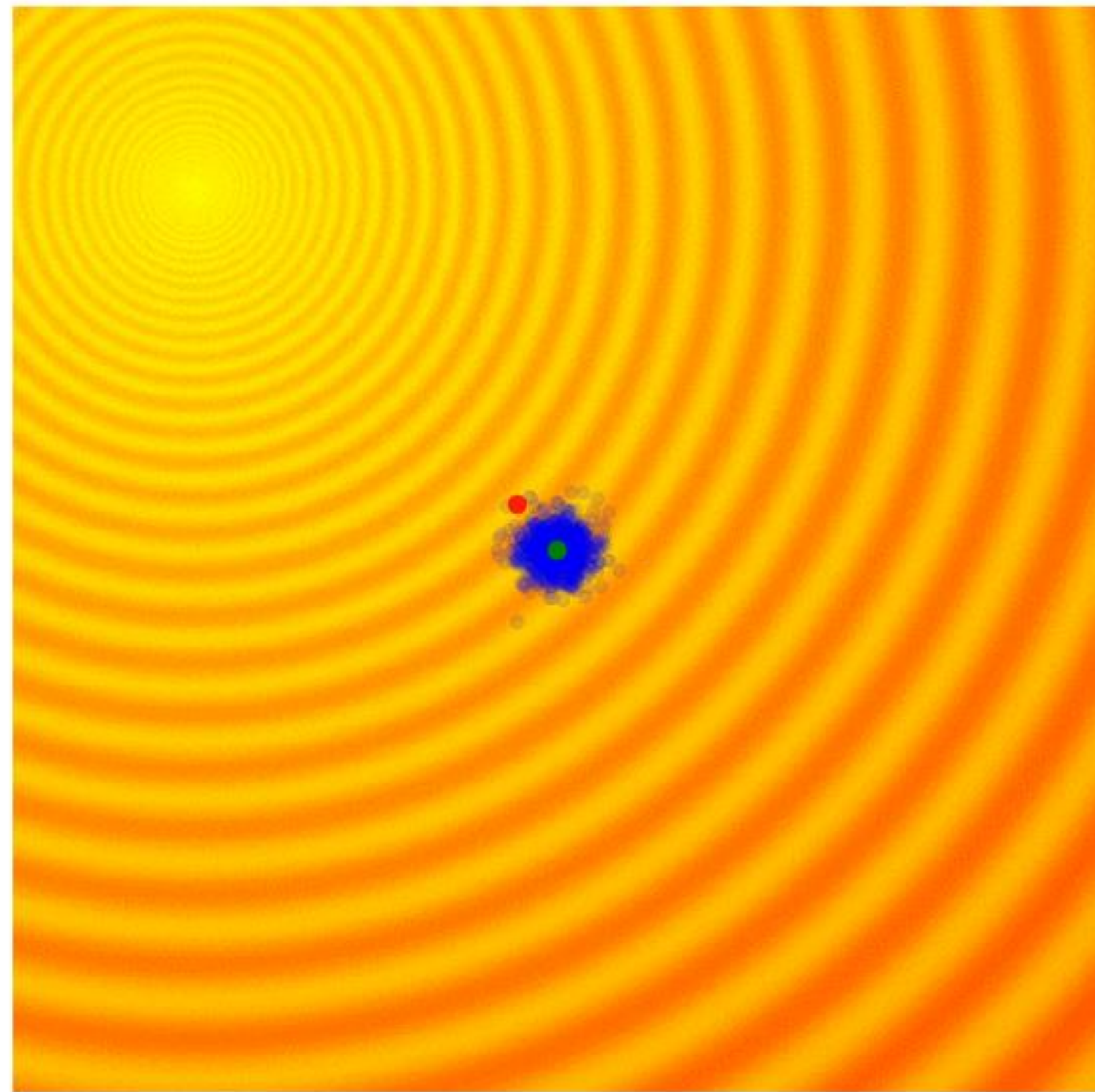
Απλές στρατηγικές εξέλιξης



[ΠΗΓΗ](#)



Προηγμένες στρατηγικές εξέλιξης: CMA-ES



[ΠΗΓΗ](#)



Neuroevolution

- Εξέλιξη των νευρωνικών δικτύων
- Εμπνευσμένο από το γεγονός ότι οι φυσικοί εγκέφαλοι είναι τα προϊόντα μιας εξελικτικής διαδικασίας
- Μπορεί να χρησιμοποιηθεί για να βελτιστοποιήσει:
 - Παράμετροι (βάρη) — αντί για gradient descent
 - Συνδεσιμότητα (τοπολογία)
 - Λειτουργίες ενεργοποίησης
 - Μαθησιακοί αλγόριθμοι (εξέλιξη της μάθησης) — πιο σχετικοί με την ενίσχυση των μαθησιακών προβλημάτων
- Ανάγκη καθορισμού της αναπαράστασης του δικτύου και των πράξεων (διασταύρωση, μετάλλαξη) που τροποποιούν την αναπαράσταση
- Μια **συμπληρωματική** προσέγγιση σε άλλους αλγόριθμους μηχανικής μάθησης!





Επόμενες Διαλέξεις

- Νευρωνικά Δίκτυα 3: Εισαγωγή στη βαθιά μάθηση

Μέρος 3: Μη εποπτευόμενη μάθηση

- Clustering



MAI4CAREU

Master programmes in Artificial
Intelligence 4 Careers in Europe



Σας ευχαριστούμε



Co-financed by the European Union
Connecting Europe Facility

This Master is run under the context of Action
No 2020-EU-IA-0087, co-financed by the EU CEF Telecom
under GA nr. INEA/CEF/ICT/A2020/2267423

