

Δείξε ότι εάν προσθέσουμε momentum στο SD, τότε πάντα υπάρχει συζεταμένος momentum που κάνει τον αλγόριθμο stable, ανεξάρτητα από το learning rate

Η steepest descent είναι

$$\Delta x_k = -\alpha \nabla F(x_k) = -\alpha g_k$$

Εάν προσθέσω momentum, έχω

$$\Delta x_k = \gamma \Delta x_{k-1} - (1-\gamma) \alpha g_k$$

Θυμίζω ότι η quadratic function είναι

$$F(x) = \frac{1}{2} x^T A x + d^T x + c$$

Η grad είναι  $\nabla F(x) = Ax + d$

Αυτοαθροισώντας έχω:

$$\Delta x_k = \gamma \Delta x_{k-1} - (1-\gamma) \alpha (Ax_k + d)$$

Αφού  $\Delta x_k = x_{k+1} - x_k$  έχω

$$x_{k+1} - x_k = \gamma (x_k - x_{k-1}) - (1-\gamma) \alpha (Ax_k + d)$$

$$\hat{n} \quad x_{k+1} = \left[ (1+\gamma)I - (1-\gamma)\alpha A \right] x_k - \gamma x_{k-1} - (1-\gamma)\alpha d$$

Ορίσω το διάνυσμα:

$$\tilde{x}_k = \begin{bmatrix} x_{k-1} \\ x_k \end{bmatrix}$$

Απα η SD γε momentum ημερα

$$\tilde{x}_{k+1} = \begin{bmatrix} 0 & I \\ -\gamma I & [(1+\gamma)I - (1-\gamma)\alpha A] \end{bmatrix} \tilde{x}_k + \begin{bmatrix} 0 \\ -(1-\gamma)\alpha d \end{bmatrix} = \underline{\underline{W\tilde{x}_k + v}}$$

Αυτο είναι linear dynamical system

Θα είναι stable, εαν οι eigenvalues του W είναι μικροτερες απο 1

Θα βρούμε τις eigenvalues σε ωδια:

$$W = \begin{bmatrix} 0 & I \\ -\gamma I & T \end{bmatrix} \quad \text{δνου } T = (1+\gamma)I - (1-\gamma)\alpha A$$

0, eigenvalues & eigenvectors του W θα ικανοποιουν

$$Wz^w = \lambda z^w \quad \text{or} \quad \begin{bmatrix} 0 & I \\ -\gamma I & T \end{bmatrix} \begin{bmatrix} z_1^w \\ z_2^w \end{bmatrix} = \lambda \begin{bmatrix} z_1^w \\ z_2^w \end{bmatrix}$$

Αυτο σημαίνει οτι

$$z_1^w = \lambda z_1^w$$

$$\text{και } -\gamma z_1^w + T z_2^w = \lambda z_2^w$$

Εστω οτι επιλεξω την  $z_2^w$  να είναι eigenvector του T γε αυτιστοιχη eigenvalue την  $\lambda^t$ . Εαν δεν είναι ορθο, θα οδηγηθω σε ατρονο (contradiction)

Αρα η προηγούμενη εξίσωση γίνεται

$$z_2^w = \lambda^w z_1^w \quad \text{και} \quad -\gamma z_1^w + \lambda^t z_2^w = \lambda^w z_2^w$$

Αντικαθιστώντας την πρώτη εξίσωση στην δεύτερη, βρίσκουμε

$$-\frac{\gamma}{\lambda^w} z_2^w + \lambda^t z_2^w = \lambda^w z_2^w \quad \Leftrightarrow \quad \left[ (\lambda^w)^2 - \lambda^t (\lambda^w) + \gamma \right] z_2^w = 0$$

Αρα για κάθε eigenvalue  $\lambda^t$  του  $T$  θα υπάρχουν δύο eigenvalues  $\lambda^w$  του  $W$  που είναι ρίζες της τετραγωνικής εξίσωσης

$$(\lambda^w)^2 - \lambda^t (\lambda^w) + \gamma = 0$$

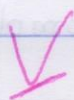
$$\text{Η λύση αυτής} \quad \lambda^w = \frac{\lambda^t \pm \sqrt{(\lambda^t)^2 - 4\gamma}}{2}$$

Για να είναι stable ο αλγόριθμος το μέτρο της κάθε eigenvalue θα πρέπει να είναι  $< 1$ . Θα δείξουμε ότι πάντα υπάρχει εύθελα μέτρο του  $\gamma$  που να ικανοποιεί αυτό.

Να σημειωθεί ότι εάν οι eigenvalues  $\lambda^w$  είναι complex, τότε το μέτρο τους είναι  $|\lambda^w| = \sqrt{\frac{(\lambda^t)^2}{4} + \frac{4\gamma - (\lambda^t)^2}{4}} = \sqrt{\gamma}$

αυτό ισχύει για πραγματικό  $\lambda^t$ . Θα δείξουμε ότι είναι όμως real

Από το  $\gamma$  είναι  $0 < \gamma < 1$ , η τιμή της eigenvalue θα πρέπει να είναι μικρότερη από 1



Απομένει να δείξουμε ότι υπάρχει εγβέλεια τιμών του  $\gamma$   
για την οποία όλες οι eigenvalues είναι complex

Για να είναι complex η  $\lambda^w$  πρέπει

$$(\lambda^t)^2 - 4\gamma < 0 \Rightarrow |\lambda^t| < 2\sqrt{\gamma}$$

As εξετάσουμε τις eigenvalues  $\lambda^t$  του  $T$

Αυτές μπορούν να γραφούν συναρτήσει των eigenvalues του  $A$

Εξω  $\{\lambda_1, \lambda_2, \dots, \lambda_n\}$  &  $\{z_1, z_2, \dots, z_n\}$  eigenpairs του Hessian

Τότε

$$Tz_i = [(1+\gamma)I - (1-\gamma)\alpha A]z_i = (1+\gamma)z_i - (1-\gamma)\alpha Az_i$$

$$= (1+\gamma)z_i - (1-\gamma)\alpha \lambda_i z_i = \left\{ (1+\gamma) - (1-\gamma)\alpha \lambda_i \right\} z_i = \lambda_i^t z_i$$

Άρα τα eigenvectors του  $T$  είναι τα ίδια με του  $A$ , και οι eigenvalues είναι

$$\lambda_i^t = \left\{ (1+\gamma) - (1-\gamma)\alpha \lambda_i \right\}$$

Η  $\lambda_i^t$  είναι real, αφού  $\gamma, \alpha$  και  $\lambda_i$  (για συμμετρικό  $A$ ) είναι real

Άρα για να είναι η  $\lambda^w$  complex, πρέπει

$$|\lambda^t| < 2\sqrt{\gamma} \quad \Leftrightarrow \quad \left| (1+\gamma) - (1-\gamma)\alpha \lambda_i \right| < 2\sqrt{\gamma}$$