

# Derin Öğrenme Deep Learning

Hazırlayan: M. Ali Akcayol  
Gazi Üniversitesi  
Bilgisayar Mühendisliği Bölümü

Bu dersin sunumları, "Deep Learning, I. Goodfellow, Y. Bengio and A. Courville, MIT Press, 2016." kitabı kullanılarak hazırlanmıştır.

## İçerik

---

- ▶ Derin modellerin eğitimi için optimizasyon
- ▶ Yapay sinir ağı optimizasyonundaki zorluklar
- ▶ Temel algoritmalar
- ▶ Parametre başlatma stratejileri
- ▶ Adaptif öğrenme oranı algoritmaları

## Derin modellerin eğitimi için optimizasyon

- ▶ Yapay sinir ağı modellerinin optimizasyonu oldukça zordur.
- ▶ Bir yapay sinir ağı probleminin çözümü bazen yüzlerce makineyle aylarca sürebilir.
- ▶ Yapay sinir ağı optimizasyonunda  $J(\theta)$  **maliyet fonksiyonunu minimum yapan  $\theta$  parametreleri bulunur.**
- ▶ Genellikle, performans ölçütü  $P$  optimize edilerek dolaylı olarak  $J$  minimum yapılır.
- ▶ Eğitim veri kümesi için maliyet fonksiyonu minimum yapılır.

$$J(\theta) = \mathbb{E}_{(\mathbf{x}, y) \sim \hat{p}_{\text{data}}} L(f(\mathbf{x}; \theta), y)$$

- ▶ Tüm veriler için maliyet fonksiyonunun minimum yapılması amaçlanır.  $m$  örnek sayısını göstermektedir.

$$J^*(\theta) = \mathbb{E}_{(\mathbf{x}, y) \sim p_{\text{data}}} L(f(\mathbf{x}; \theta), y)$$

$$\mathbb{E}_{\mathbf{x}, y \sim \hat{p}_{\text{data}}(\mathbf{x}, y)} [L(f(\mathbf{x}; \theta), y)] = \frac{1}{m} \sum_{i=1}^m L(f(\mathbf{x}^{(i)}; \theta), y^{(i)})$$

3

## Derin modellerin eğitimi için optimizasyon

- ▶ Eğitim veri kümesinin tamamını kullanan optimizasyon algoritmaları **batch** veya **deterministik** gradyan metot olarak adlandırılır.
- ▶ Sadece bir tane örnek veriyi kullanan optimizasyon algoritmaları **online** veya **stokastik** metot olarak adlandırılır.
- ▶ Derin öğrenme algoritmalarının çoğu birden çok ancak tüm eğitim veri kümesinden daha az örnek veri kullanır.
- ▶ Bu tür algoritmalar **minibatch** veya **minibatch stokastik** metot olarak adlandırılır.

4

## Derin modellerin eğitimi için optimizasyon

- ▶  $J(\mathbf{X})$  maliyet fonksiyonunu minimum yapacak güncelleme değeri minibatch  $X$  için hesaplanır.
- ▶ Genelleme hatası ve gradyan aşağıdaki şekilde ifade edilir:

$$J^*(\boldsymbol{\theta}) = \sum_{\mathbf{x}} \sum_y p_{\text{data}}(\mathbf{x}, y) L(f(\mathbf{x}; \boldsymbol{\theta}), y)$$

$$\mathbf{g} = \nabla_{\boldsymbol{\theta}} J^*(\boldsymbol{\theta}) = \sum_{\mathbf{x}} \sum_y p_{\text{data}}(\mathbf{x}, y) \nabla_{\boldsymbol{\theta}} L(f(\mathbf{x}; \boldsymbol{\theta}), y)$$

$$\hat{\mathbf{g}} = \frac{1}{m} \nabla_{\boldsymbol{\theta}} \sum_i L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), y^{(i)})$$

5

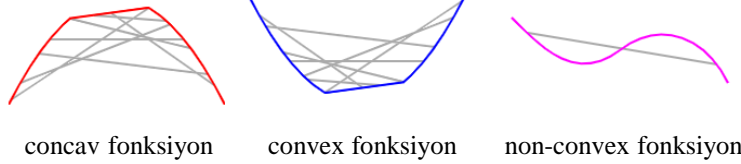
## İçerik

- ▶ Derin modellerin eğitimi için optimizasyon
- ▶ Yapay sinir ağı optimizasyonundaki zorluklar
- ▶ Temel algoritmalar
- ▶ Parametre başlatma stratejileri
- ▶ Adaptif öğrenme oranı algoritmaları

6

## Yapay sinir ağı optimizasyonundaki zorluklar

- ▶ Optimizasyon genellikle çok zordur ve uzun sürer.
- ▶ Makine öğrenmesinde **amaç fonksiyonu ve kısıtlar** optimizasyon problemini **convex yapacak şekilde tasarlanır.**

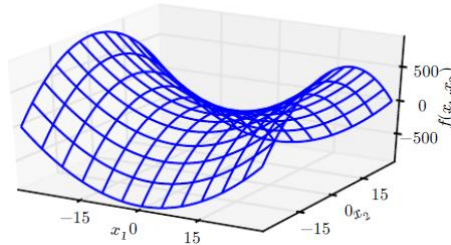


- ▶ **Genellikle,** yapay sinir ağı eğitiminde **non-convex durumla karşılaşılır.**
- ▶ Convex optimizasyon probleminde lokal minimum aynı zamanda global minimumdur.
- ▶ Non-convex fonksiyonlarda birden fazla lokal minimum olabilir.
- ▶ Derin öğrenme modellerinde çok sayıda lokal minimum olabilir.

7

## Yapay sinir ağı optimizasyonundaki zorluklar

- ▶ Yüksek boyutlu non-convex fonksiyonlar sıfır gradyana sahip olan **saddle point'e (eyer noktası) sahip olabilir.**
- ▶ Bir eyer noktasında Hessian matris **hem pozitif hem de negatif** eigenvalue'ya sahiptir (hem maksimum hem de minimum nokta vardır.).

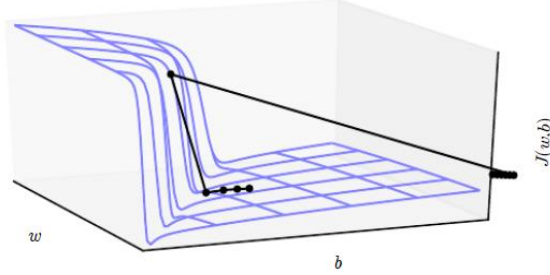


- ▶ Düşük boyutlu uzaylarda **local minimum yaygındır.**
- ▶ Yüksek boyutlu uzaylarda **saddle point yaygındır.**
- ▶  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  fonksiyonu için saddle point sayısının local minimum sayısına oranının üstel artması beklenir.

8

## Yapay sinir ađı optimizasyonundaki zorluklar

- ▶ Birçok yapay sinir ađı uçurumlara benzer **aşırı dik bölgelere sahiptir.**
- ▶ Bunun nedeni **çok sayıda büyük değere sahip ađırlığın birbiriyle çarpılmasıdır.**



- ▶ Bu durumda **gradyan güncelleme parametreleri çok büyük oranda deđiştir.**
- ▶ Klasik **gradient descent algoritması** parametreler için **büyük bir deđişim değeri hesaplar.**

9

## İçerik

- ▶ Derin modellerin eğitimi için optimizasyon
- ▶ Yapay sinir ađı optimizasyonundaki zorluklar
- ▶ **Temel algoritmalar**
- ▶ Parametre başlatma stratejileri
- ▶ Adaptif öğrenme oranı algoritmaları

10

## Temel algoritmalar

### Stochastic gradient descent

- ▶ Derin öğrenmede en çok kullanılan optimizasyon algoritmasıdır.

Algorithm 8.1 Stochastic gradient descent (SGD) update at training iteration  $k$

Require: Learning rate  $\epsilon_k$ .

Require: Initial parameter  $\theta$

while stopping criterion not met do

    Sample a minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with corresponding targets  $\mathbf{y}^{(i)}$ .

    Compute gradient estimate:  $\hat{\mathbf{g}} \leftarrow +\frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

    Apply update:  $\theta \leftarrow \theta - \epsilon \hat{\mathbf{g}}$

end while

- ▶ SGD algoritmasının en önemli parametresi öğrenme oranıdır.
- ▶ **Öğrenme oranı iterasyon sayısı arttıkça azaltılır** ( $\epsilon_k$ ).

$$\epsilon_k = (1 - \alpha)\epsilon_0 + \alpha\epsilon_{\tau} \quad \alpha = \frac{k}{\tau}$$

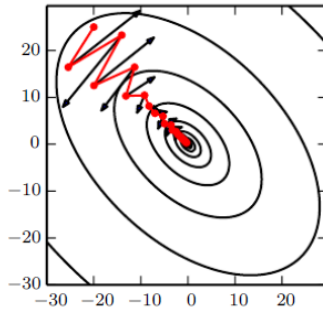
- ▶  $\epsilon$  iterasyon sonra genellikle  $\tau$  sabit bırakılır.

11

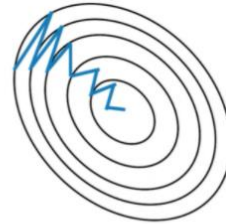
## Temel algoritmalar

### Momentum

- ▶ **Stochastic gradient descent** çok popüler bir optimizasyon algoritması olmasına rağmen **bazen oldukça yavaştır**.
- ▶ **Momentum** metodu **öğrenmeyi hızlandırmak için** kullanılır.



Stochastic Gradient  
Descent **without**  
Momentum



Stochastic Gradient  
Descent **with**  
Momentum

$$\mathbf{v} \leftarrow \alpha \mathbf{v} - \epsilon \nabla_{\theta} \left( \frac{1}{m} \sum_{i=1}^m L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)}) \right) \quad \alpha \in [0, 1]$$

$$\theta \leftarrow \theta + \mathbf{v}$$

12

## Temel algoritmalar

### Momentum

- ▶ Momentum ile stochastic gradient descent algoritması

---

Algorithm 8.2 Stochastic gradient descent (SGD) with momentum

---

Require: Learning rate  $\epsilon$ , momentum parameter  $\alpha$ .

Require: Initial parameter  $\theta$ , initial velocity  $v$ .

while stopping criterion not met do

    Sample a minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with corresponding targets  $\mathbf{y}^{(i)}$ .

    Compute gradient estimate:  $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

    Compute velocity update:  $\mathbf{v} \leftarrow \alpha \mathbf{v} - \epsilon \mathbf{g}$

    Apply update:  $\theta \leftarrow \theta + \mathbf{v}$

end while

---

13

## Temel algoritmalar

### Nesterov momentum

- ▶ Nesterov momentum'da parametreler **önce mevcut değişim hızıyla değiştirilir.**
- ▶ Ardından standart momentum'daki gibi tekrar değiştirilir.

$$\mathbf{v} \leftarrow \alpha \mathbf{v} - \epsilon \nabla_{\theta} \left[ \frac{1}{m} \sum_{i=1}^m L(\mathbf{f}(\mathbf{x}^{(i)}; \theta + \alpha \mathbf{v}), \mathbf{y}^{(i)}) \right]$$
$$\theta \leftarrow \theta + \mathbf{v}$$

---

Algorithm 8.3 Stochastic gradient descent (SGD) with Nesterov momentum

---

Require: Learning rate  $\epsilon$ , momentum parameter  $\alpha$ .

Require: Initial parameter  $\theta$ , initial velocity  $v$ .

while stopping criterion not met do

    Sample a minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with corresponding labels  $\mathbf{y}^{(i)}$ .

    Apply interim update:  $\tilde{\theta} \leftarrow \theta + \alpha \mathbf{v}$

    Compute gradient (at interim point):  $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\tilde{\theta}} \sum_i L(f(\mathbf{x}^{(i)}; \tilde{\theta}), \mathbf{y}^{(i)})$

    Compute velocity update:  $\mathbf{v} \leftarrow \alpha \mathbf{v} - \epsilon \mathbf{g}$

    Apply update:  $\theta \leftarrow \theta + \mathbf{v}$

end while

---

14

## İçerik

---

- ▶ Derin modellerin eğitimi için optimizasyon
- ▶ Yapay sinir ağı optimizasyonundaki zorluklar
- ▶ Temel algoritmalar
- ▶ **Parametre başlatma stratejileri**
- ▶ Adaptif öğrenme oranı algoritmaları

15

## Parametre başlatma stratejileri

---

- ▶ Derin öğrenme **eğitim algoritmaları iteratif çalışır.**
- ▶ İterasyonların başlaması için **kullanıcının bir başlangıç noktası belirlemesi gerekir.**
- ▶ Derin öğrenme algoritmalarının çoğu **başlangıç noktası seçiminden etkilenir.**
- ▶ Seçilen başlangıç noktası algoritmanın hiçbir zaman istenen sonuca yakınsamamasına neden olabilir.
- ▶ Seçilen başlangıç noktası algoritmanın istenen sonuca yakınsasa bile çok uzun süre gerekmesine neden olabilir.
- ▶ Farklı birimlerin (gizli katmanlar) ağırlıkları farklı şekilde başlatılmalıdır. Aynı değere sahip olmamalıdır.
- ▶ Her birimin farklı değerler alması için **parametreler rastgele başlatılır.**

16



## Parametre başlatma stratejileri

- ▶ Yüksek boyutlu uzayda **yüksek entropi dağılımına sahip rastgele başlatmanın hesaplama maliyeti düşüktür.**
- ▶ Genellikle, **bias değerleri sabit alınır** ve **ağırlık değerleri rastgele başlatılır.**
- ▶ Tüm ağırlık değerleri genellikle **Gauss veya uniform dağılımla** rastgele atanır.
- ▶ Büyük başlangıç değerleri, doğrusal fonksiyona sahip düğümlerde verinin hem ileri geçişte hem de geri yayılımda kaybolmasına engel olur.
- ▶ Çok büyük başlangıç değerleri ise değerlerde hem ileri geçişte hem de geri yayılımda patlayan değerlere neden olur.
- ▶ **Recurrent neural network'lerde büyük ağırlık değerleri kaos'a neden olur.**

17

## Parametre başlatma stratejileri

### Xavier initialization

- ▶  $m$  tane giriş için aşağıdaki gibi ifade edilir.

$$U\left(-\frac{1}{\sqrt{m}}, \frac{1}{\sqrt{m}}\right)$$

- ▶ Burada,  $U$  uniform dağılımı gösterir.

### Normalized initialization

- ▶ Glorot ve Bengio'nun önerdiği normalize edilmiş başlatma aşağıdaki şekilde ifade edilir.

$$W_{i,j} \sim U\left(-\sqrt{\frac{6}{m+n}}, \sqrt{\frac{6}{m+n}}\right)$$

- ▶ Burada,  $m$  giriş sayısını,  $n$  çıkış sayısını ve  $U$  uniform dağılımı gösterir.

18

## İçerik

---

- ▶ Derin modellerin eğitimi için optimizasyon
- ▶ Yapay sinir ağı optimizasyonundaki zorluklar
- ▶ Temel algoritmalar
- ▶ Parametre başlatma stratejileri
- ▶ **Adaptif öğrenme oranı algoritmaları**

19

## Adaptif öğrenme oranı algoritmaları

---

### **Delta-bar-delta algoritması**

- ▶ **Öğrenme oranı her bir parametre için** kısmi türev değerine göre **değiştirilir**.
- ▶ Eğer kayıp fonksiyonunun model parametrelerine göre **kısmi türevi, parametre işaretini aynı yönde bırakıyorsa**, öğrenme oranı **artırılır**.
- ▶ Eğer kayıp fonksiyonunun kısmi türevi parametre işaretini **değiştiriyorsa**, öğrenme oranı **azaltılır**.
- ▶ Bu kural full batch optimizasyonda da kullanılabilir.

20

## Adaptif öğrenme oranı algoritmaları

### AdaGrad algoritması

- ▶ AdaGrad algoritması ile öğrenme oranı her bir parametre için **tüm geçmiş değerlerinin toplamının karekökü ile orantılı olarak değiştirilir.**
- ▶ **Kısmi türev değeri büyük olan** bir parametre için öğrenme oranı **hızlı azalır, küçük olan için azalma yavaştır.**
- ▶ Başlangıçtan itibaren gradyan değerlerinin toplanarak alınması erken ve fazla düşüğe neden olabilir.
- ▶ AdaGrad algoritması bazı derin öğrenme modellerinde başarılı olmuştur ancak **hepsinde başarılı değildir.**

21

## Adaptif öğrenme oranı algoritmaları

### AdaGrad algoritması

Algorithm 8.4 The AdaGrad algorithm

Require: Global learning rate  $\epsilon$

Require: Initial parameter  $\theta$

Require: Small constant  $\delta$ , perhaps  $10^{-7}$ , for numerical stability

Initialize gradient accumulation variable  $\mathbf{r} = \mathbf{0}$

while stopping criterion not met do

    Sample a minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with corresponding targets  $\mathbf{y}^{(i)}$ .

    Compute gradient:  $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

    Accumulate squared gradient:  $\mathbf{r} \leftarrow \mathbf{r} + \mathbf{g} \odot \mathbf{g}$

    Compute update:  $\Delta \theta \leftarrow -\frac{\epsilon}{\delta + \sqrt{\mathbf{r}}} \odot \mathbf{g}$ . (Division and square root applied element-wise)

    Apply update:  $\theta \leftarrow \theta + \Delta \theta$

end while

22

## Adaptif öğrenme oranı algoritmaları

### RMSPProp algoritması

- ▶ RMSPProp algoritması AdaGrad algoritmasındaki gradyan toplama işlemini değiştirmiştir.
- ▶ AdaGrad, convex fonksiyonlarda hızlı yakınsama için tasarlanmıştır.
- ▶ AdaGrad, non-convex fonksiyonlarda uygulandığında lokal convex bir bölgeye ulaşabilmektedir.
- ▶ RMSPProp algoritması **uzun süreli geçmişli atarak convex bir bölgeye hızlı yakınsama yapabilir.**
- ▶ RMSPProp algoritması bir **hiperparametre ( $\rho$ ) kullanarak önceki değeri ile gradyan karesi arasında ağırlıklandırma yapar.**
- ▶ RMSPProp, derin sinirsel ağlarda etkin bir optimizasyon algoritmasıdır.

23

## Adaptif öğrenme oranı algoritmaları

### RMSPProp algoritması

Algorithm 8.5 The RMSPProp algorithm

Require: Global learning rate  $\epsilon$ , decay rate  $\rho$ .

Require: Initial parameter  $\theta$

Require: Small constant  $\delta$ , usually  $10^{-6}$ , used to stabilize division by small numbers.

Initialize accumulation variables  $r = 0$

while stopping criterion not met do

    Sample a minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with corresponding targets  $\mathbf{y}^{(i)}$ .

    Compute gradient:  $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

    Accumulate squared gradient:  $\mathbf{r} \leftarrow \rho \mathbf{r} + (1 - \rho) \mathbf{g} \odot \mathbf{g}$

    Compute parameter update:  $\Delta \theta = -\frac{\epsilon}{\sqrt{\delta + \mathbf{r}}} \odot \mathbf{g}$ . ( $\frac{1}{\sqrt{\delta + \mathbf{r}}}$  applied element-wise)

    Apply update:  $\theta \leftarrow \theta + \Delta \theta$

end while

24

## Adaptif öğrenme oranı algoritmaları

### RMSProp algoritması ve Nesterov momentum

Algorithm 8.6 RMSProp algorithm with Nesterov momentum

Require: Global learning rate  $\epsilon$ , decay rate  $\rho$ , momentum coefficient  $\alpha$ .

Require: Initial parameter  $\theta$ , initial velocity  $v$ .

Initialize accumulation variable  $r = \mathbf{0}$

while stopping criterion not met do

Sample a minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with corresponding targets  $\mathbf{y}^{(i)}$ .

Compute interim update:  $\tilde{\theta} \leftarrow \theta + \alpha v$

Compute gradient:  $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\tilde{\theta}} \sum_i L(f(\mathbf{x}^{(i)}; \tilde{\theta}), \mathbf{y}^{(i)})$

Accumulate gradient:  $\mathbf{r} \leftarrow \rho \mathbf{r} + (1 - \rho) \mathbf{g} \odot \mathbf{g}$

Compute velocity update:  $\mathbf{v} \leftarrow \alpha v - \frac{\epsilon}{\sqrt{\mathbf{r}}} \odot \mathbf{g}$ . ( $\frac{1}{\sqrt{\mathbf{r}}}$  applied element-wise)

Apply update:  $\theta \leftarrow \theta + v$

end while

25

## Adaptif öğrenme oranı algoritmaları

### Adam (adaptive moments) algoritması

► Adam, birinci ve ikinci derece gradyan değerlerini kullanır.

Algorithm 8.7 The Adam algorithm

Require: Step size  $\epsilon$  (Suggested default: 0.001)

Require: Exponential decay rates for moment estimates,  $\rho_1$  and  $\rho_2$  in  $[0, 1)$ . (Suggested defaults: 0.9 and 0.999 respectively)

Require: Small constant  $\delta$  used for numerical stabilization. (Suggested default:  $10^{-8}$ )

Require: Initial parameters  $\theta$

Initialize 1st and 2nd moment variables  $\mathbf{s} = \mathbf{0}$ ,  $\mathbf{r} = \mathbf{0}$

Initialize time step  $t = 0$

while stopping criterion not met do

Sample a minibatch of  $m$  examples from the training set  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  with corresponding targets  $\mathbf{y}^{(i)}$ .

Compute gradient:  $\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

$t \leftarrow t + 1$

Update biased first moment estimate:  $\mathbf{s} \leftarrow \rho_1 \mathbf{s} + (1 - \rho_1) \mathbf{g}$

Update biased second moment estimate:  $\mathbf{r} \leftarrow \rho_2 \mathbf{r} + (1 - \rho_2) \mathbf{g} \odot \mathbf{g}$

Correct bias in first moment:  $\hat{\mathbf{s}} \leftarrow \frac{\mathbf{s}}{1 - \rho_1^t}$

Correct bias in second moment:  $\hat{\mathbf{r}} \leftarrow \frac{\mathbf{r}}{1 - \rho_2^t}$

Compute update:  $\Delta \theta = -\epsilon \frac{\hat{\mathbf{s}}}{\sqrt{\hat{\mathbf{r}} + \delta}}$  (operations applied element-wise)

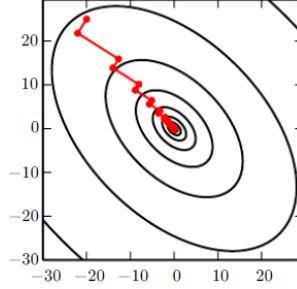
Apply update:  $\theta \leftarrow \theta + \Delta \theta$

end while

## Adaptif öğrenme oranı algoritmaları

### Conjugate gradients

- ▶ Conjugate gradients, **arama yönündeki her doğrunun bir öncekine dik olmasını garanti eder.**



- ▶ İterasyon  $t$  için  $d_t$  aşağıdaki şekilde hesaplanır ( $\beta_t$  sabit katsayı):

$$d_t = \nabla_{\theta} J(\theta) + \beta_t d_{t-1}$$

- ▶ Eğer  $d_t^T H d_{t-1} = 0$  ise,  $d_t$  ve  $d_{t-1}$  conjugate olarak ifade edilir.

$$H \leftarrow \frac{1}{m} \nabla_{\theta}^2 \sum_i L(f(x^{(i)}; \theta), y^{(i)})$$

27

## Adaptif öğrenme oranı algoritmaları

### Conjugate gradients

- ▶  $\beta_t$  katsayısının belirlenmesi için iki popüler yöntem vardır:

Fletcher-Reeves:

$$\beta_t = \frac{\nabla_{\theta} J(\theta_t)^T \nabla_{\theta} J(\theta_t)}{\nabla_{\theta} J(\theta_{t-1})^T \nabla_{\theta} J(\theta_{t-1})}$$

Polak-Ribière:

$$\beta_t = \frac{(\nabla_{\theta} J(\theta_t) - \nabla_{\theta} J(\theta_{t-1}))^T \nabla_{\theta} J(\theta_t)}{\nabla_{\theta} J(\theta_{t-1})^T \nabla_{\theta} J(\theta_{t-1})}$$

- ▶  $k$ -boyutlu bir parametre uzayında conjugate gradients metodu en çok  $k$  çizgi aramasıyla minimuma ulaşır.

28

## Adaptif öğrenme oranı algoritmaları

### Conjugate gradients

Algorithm 8.9 The conjugate gradient method

Require: Initial parameters  $\theta_0$

Require: Training set of  $m$  examples

Initialize  $\rho_0 = \mathbf{0}$

Initialize  $g_0 = \mathbf{0}$

Initialize  $t = 1$

while stopping criterion not met do

  Initialize the gradient  $g_t = \mathbf{0}$

  Compute gradient:  $g_t \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

  Compute  $\beta_t = \frac{(g_t - g_{t-1})^\top g_t}{g_{t-1}^\top g_{t-1}}$  (Polak-Ribière)

  (Nonlinear conjugate gradient: optionally reset  $\beta_t$  to zero, for example if  $t$  is a multiple of some constant  $k$ , such as  $k = 5$ )

  Compute search direction:  $\rho_t = -g_t + \beta_t \rho_{t-1}$

  Perform line search to find:  $\epsilon^* = \operatorname{argmin}_{\epsilon} \frac{1}{m} \sum_{i=1}^m L(f(\mathbf{x}^{(i)}; \theta_t + \epsilon \rho_t), \mathbf{y}^{(i)})$

  (On a truly quadratic cost function, analytically solve for  $\epsilon^*$  rather than explicitly searching for it)

  Apply update:  $\theta_{t+1} = \theta_t + \epsilon^* \rho_t$

$t \leftarrow t + 1$

end while