

# Supervised Learning

Yaya Wihardi, S.Kom., M.Kom.

Email: [yayawihardi@upi.edu](mailto:yayawihardi@upi.edu)

**Department of Computer Science Education**

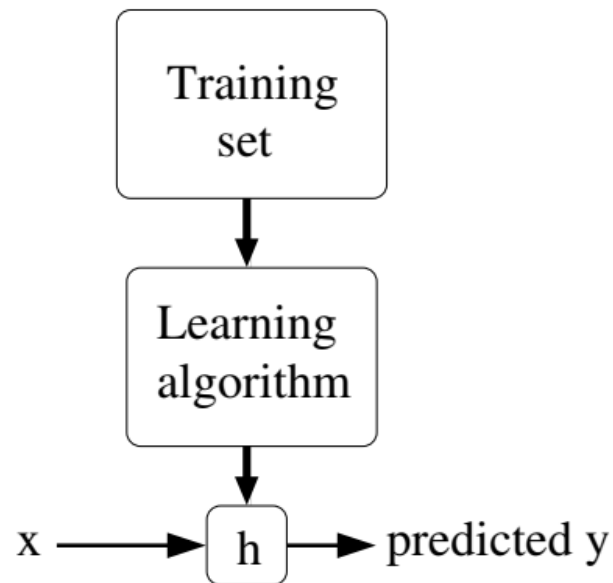
Universitas Pendidikan Indonesia

# Outline

- Apa itu Supervised Learning?
- Task dalam Supervised Learning
- Notasi yang sering digunakan
- Regresi Linear
- Least Mean Square
- Gradient Descent

# Apa itu Supervised Learning?

- Proses pembelajaran mesin yang mencoba membuat fungsi pemetaan dari input menjadi output berdasarkan data training



# Task dalam Supervised Learning

- Regression → memprediksi nilai continue, seperti harga rumah, harga saham, dll
- Classification → Memprediksi nilai diskrit, seperti kategori, kelas, kelompok, dll



**Apakah ini foto Susi, Steak, atau Pizza?**



# Your Turn

- Anda diberikan history data penjualan item-item produk dalam 3 tahun terakhir, prediksi berapa jumlah penjualan setiap item produk dalam 3 bulan kedepan.
- Anda diberikan data akun dan riwayat aktifitasnya dalam sebuah sistem pemerintahan. Prediksi apakah akun-akun tersebut diretas atau tidak.
- Diberikan data email yang telah dilabeli spam atau bukan spam. Buatlah spam filter.

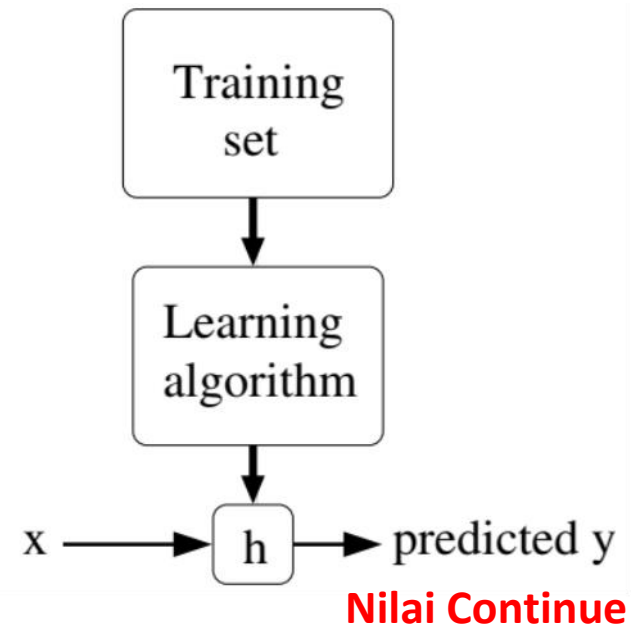
# Notasi

Notasi	Keterangan
$x^{(i)}$	Variabel Input/Fitur
$y^{(i)}$	Variabel Output/Target/ Variabel yang akan diprediksi
$(x^{(i)}, y^{(i)})$	Pasangan Training Example/Instance/Data Point
$\{(x^{(i)}, y^{(i)}); i = 1, \dots, n\}$	Dataset/Himpunan Data yg digunakan untuk membangun model
$\dots^{(i)}$	(i) Menyatakan indeks dari data, <b>bukan pangkat</b>
$R^n$	Ruang vector berdimensi n

# Regresi: Studi Kasus Real Data

- Diberikan data luas living area (feet<sup>2</sup>), jml bedrooms, dan harganya (USD) sbb:

Living area (feet <sup>2</sup> )	#bedrooms	Price (1000\$s)
2104	3	400
1600	3	330
2400	3	369
1416	2	232
3000	4	540
⋮	⋮	⋮



- Buat model ML yg dapat memprediksi harga rumah berdasarkan variable input luas living area dan jml bedrooms

# Regresi: Studi Kasus Real Data

- Diberikan data luas living area (feet<sup>2</sup>), jml bedrooms, dan harganya (USD) sbb:

Living area (feet <sup>2</sup> )	#bedrooms	Price (1000\$s)
2104	3	400
1600	3	330
2400	3	369
1416	2	232
3000	4	540
⋮	⋮	⋮

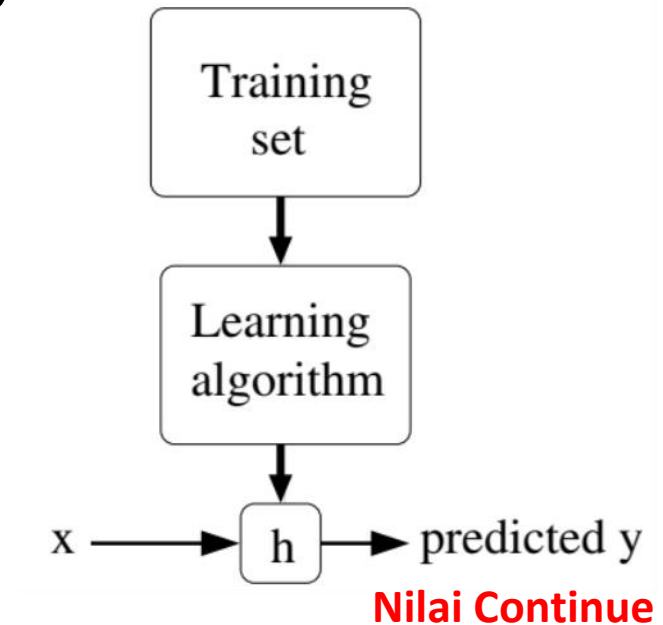
- Utk kasus ini,  $x$  merupakan vector berdimensi 2 ( $R^2$ )
- $x_1^{(i)}$  menyatakan living area,  $x_2^{(i)}$  menyatakan jml bedrooms



# Regresi Linear: Studi Kasus Sederhana

- Diberikan data pasangan  $x$  dan  $y$  sbb:

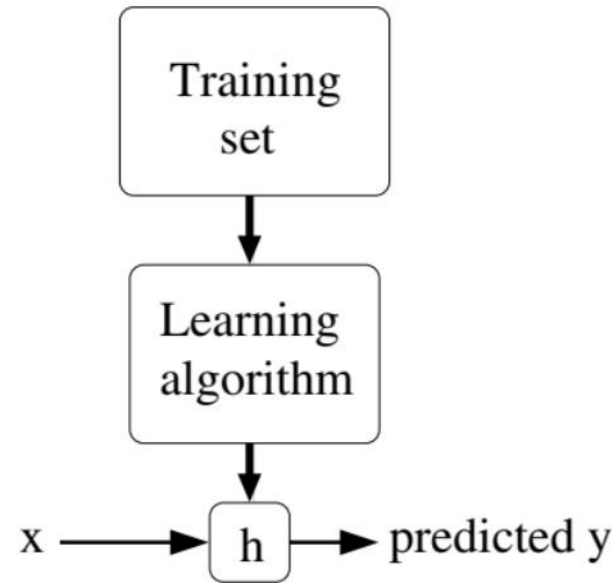
$x$	$y$
2	8
5	17
4	13
7	24
8	26
.	.
.	.
.	.



- Buat model ML yg dapat memprediksi nilai  $y$  berdasarkan input  $x$

# Caranya?

x	y
2	8
5	17
4	13
7	24
8	26
.	.
.	.
.	.



- Cari fungsi hipotesis ( $h$ ) yang dapat memetakan  $x$  ke  $y$
- Ingat prinsip occam's razor: mulailah dari yang sederhana
- Sebagai contoh, bisa digunakan fungsi persamaan linear

# Model Hipotesis (h)

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

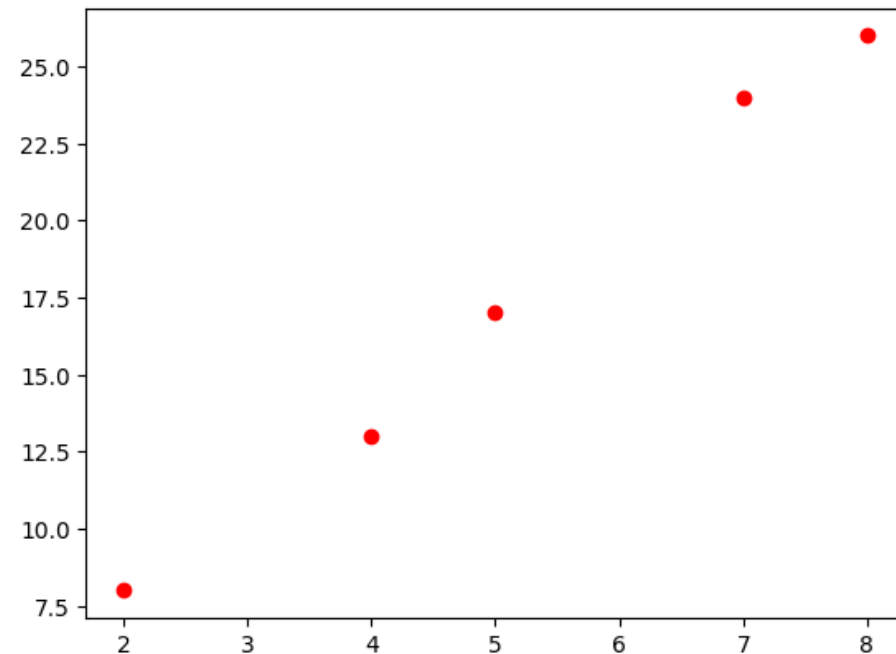
$$h_w(x) = w_0 + w_1 x$$

$\theta_i$  : Parameter model

$w_i$  : Parameter model

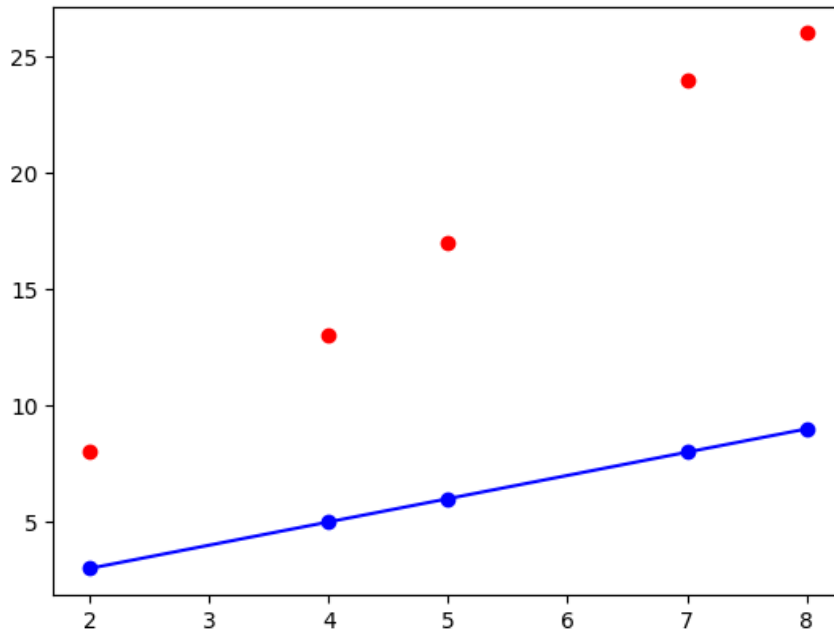
x	y
2	8
5	17
4	13
7	24
8	26
.	.

**Bagaimana memilih parameter  $w_i/\theta_i$  yang tepat?**



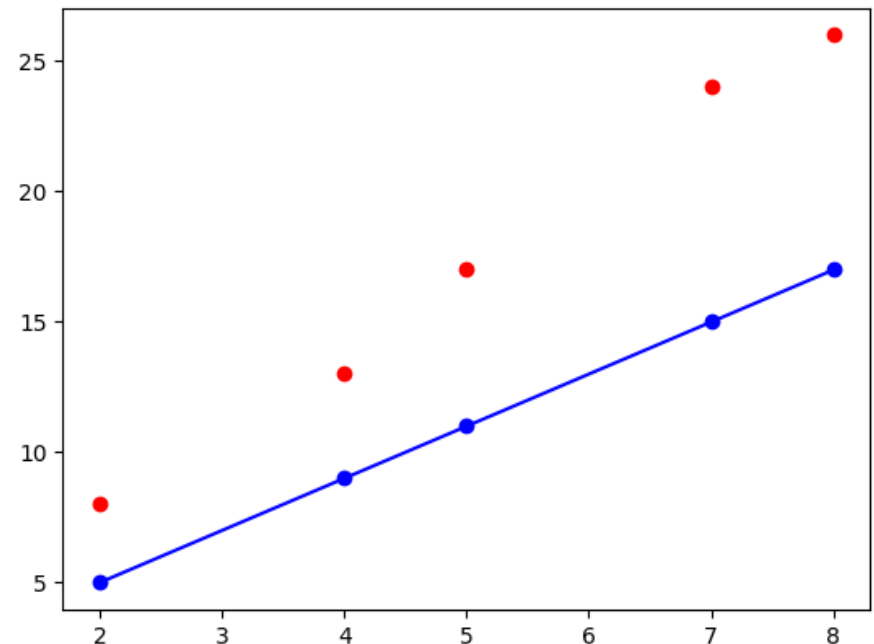
# Model Hipotesis

$$h_w(x) = w_0 + w_1 x$$



$$w_0 = 1$$

$$w_1 = 1$$

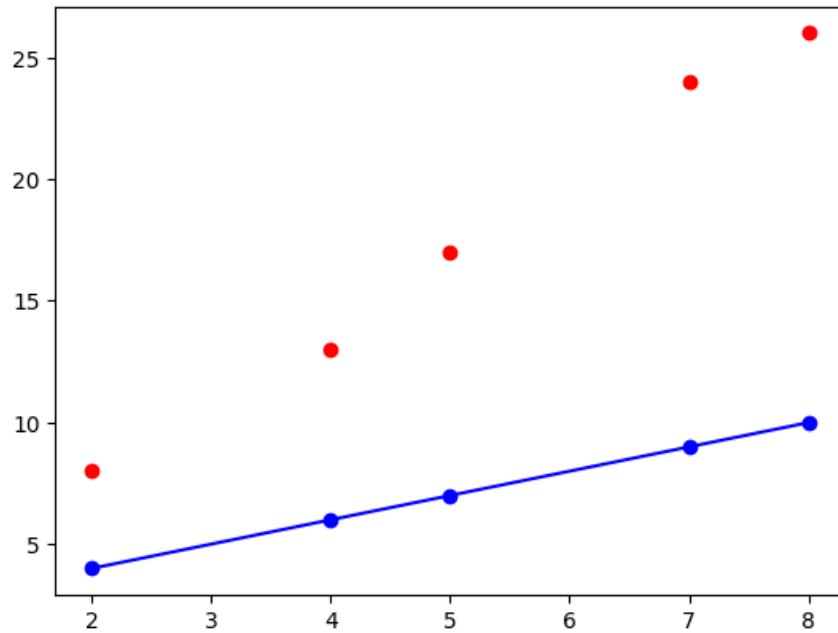


$$w_0 = 1$$

$$w_1 = 2$$

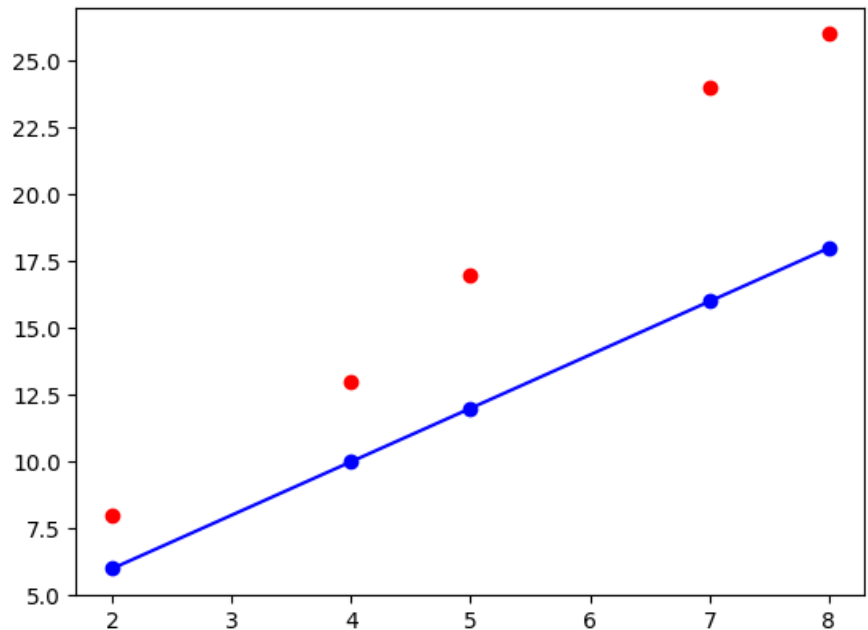
# Model Hipotesis

$$h_w(x) = w_0 + w_1 x$$



$$w_0 = 2$$

$$w_1 = 1$$



$$w_0 = 2$$

$$w_1 = 2$$

# Memilih Model Hipotesis (h)

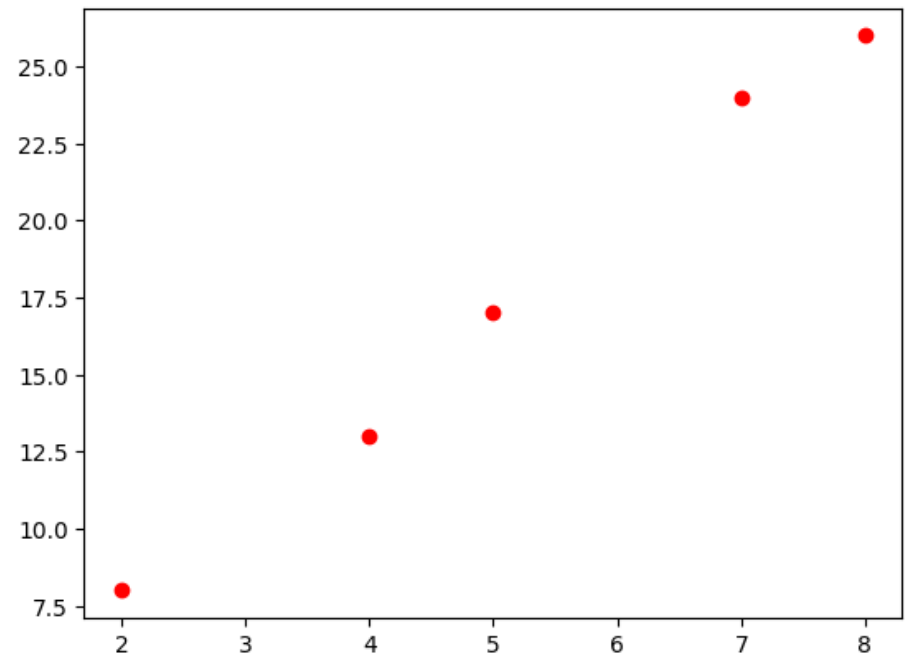
$$h_w(x) = w_0 + w_1x$$

$w_i$  : Parameter model

**Bagaimana memilih parameter  $w_i$  yang tepat?**

**IDE:** Pilihlah  $w_0, w_1$  yang membuat  $h_w(x)$  paling mendekati  $y$  untuk semua pasangan  $(x, y)$

x	y
2	8
5	17
4	13
7	24
8	26
.	.



# Cost Function: Least Mean Square (LMS)

Hipotesis:

$$h_w(x) = w_0 + w_1x$$

Parameter:

$$w_0, w_1$$

Cost Function:

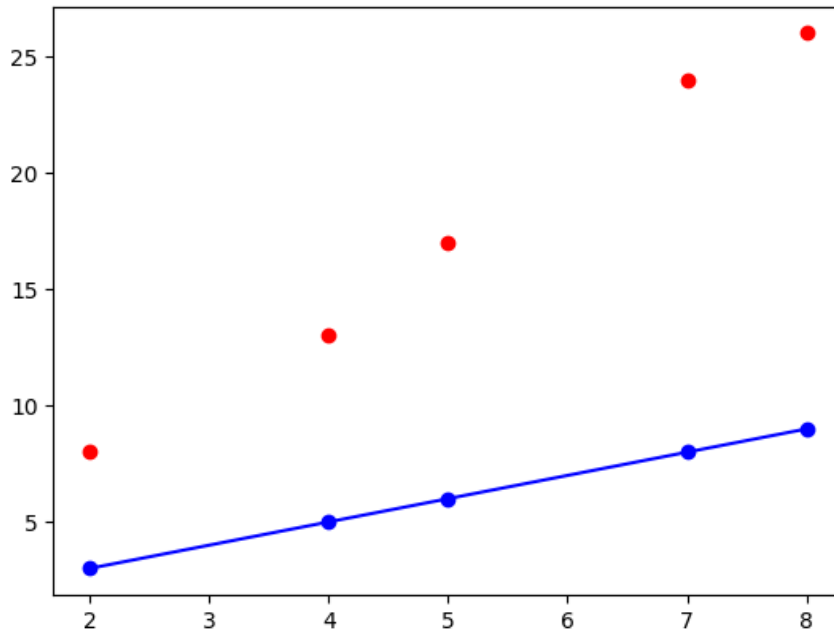
$$J(w_0, w_1) = \frac{1}{2m} \sum_{i=1}^m (h_w(x^{(i)}) - y^{(i)})^2$$

Goal:

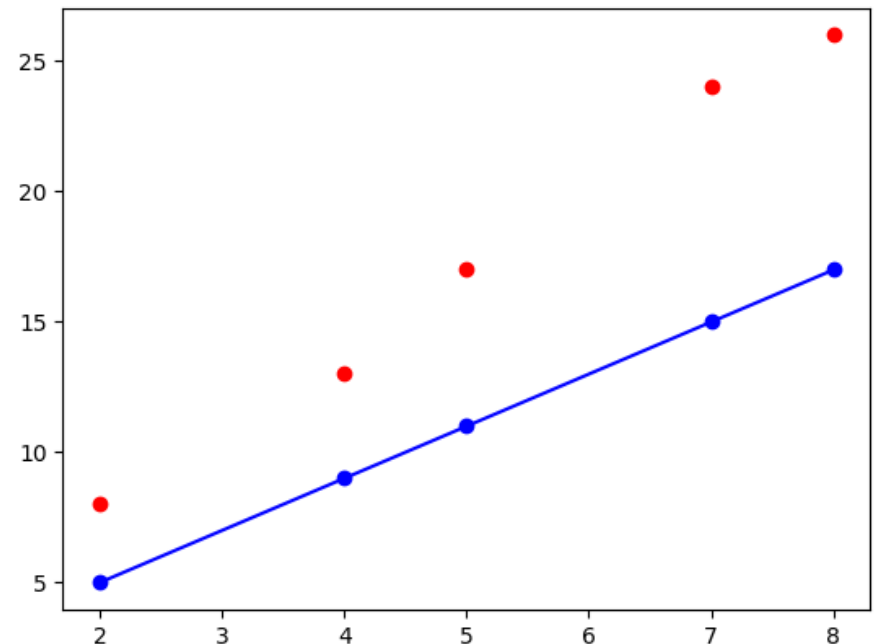
$$\underset{w_0, w_1}{\text{minimize}} \quad J(w_0, w_1)$$

# Model Hipotesis

$$h_w(x) = w_0 + w_1 x$$



$$w_0 = 1 \quad J = 75.5$$
$$w_1 = 1$$

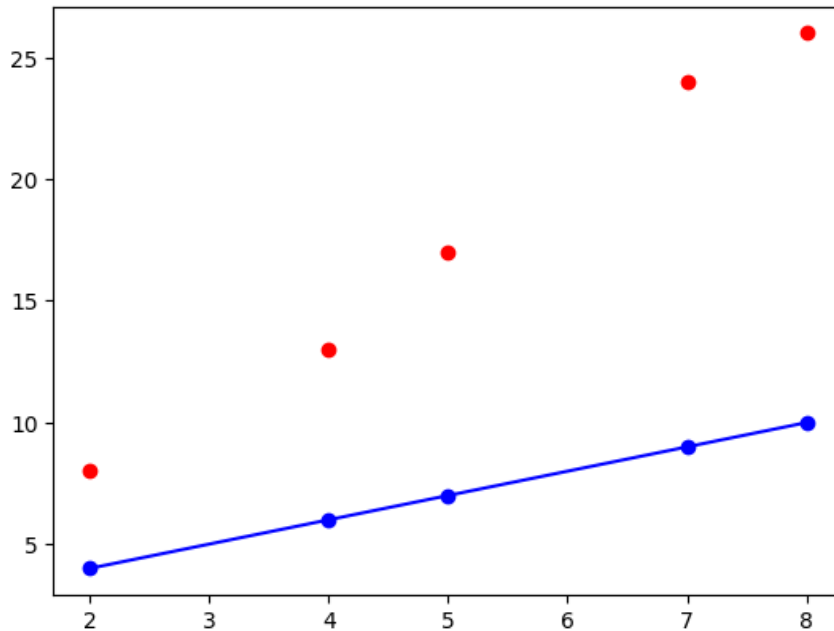


$$w_0 = 1 \quad J = 22.3$$
$$w_1 = 2$$

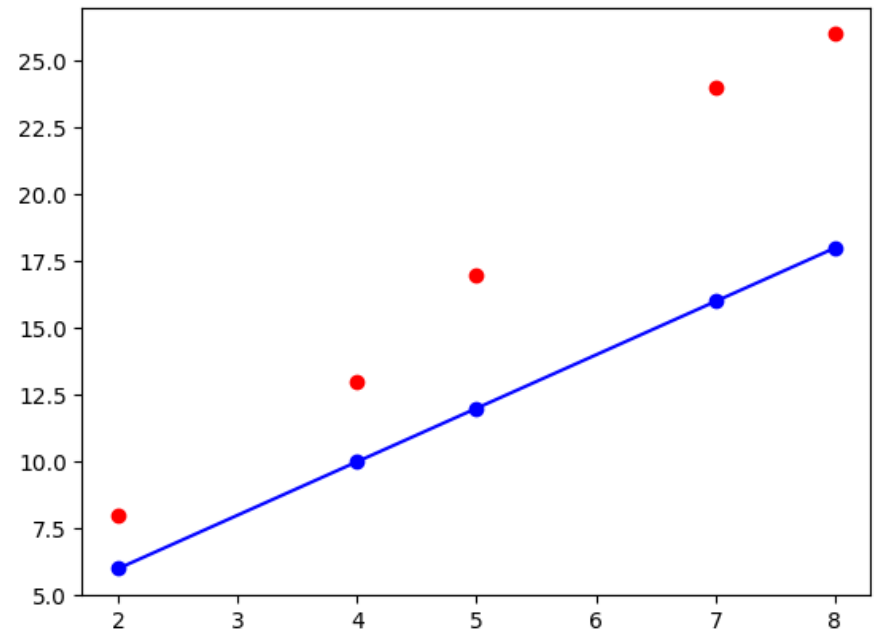


# Model Hipotesis

$$h_w(x) = w_0 + w_1 x$$



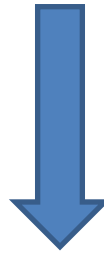
$$w_0 = 2 \quad J = 64.6$$
$$w_1 = 1$$



$$w_0 = 2 \quad J = 16.6$$
$$w_1 = 2$$

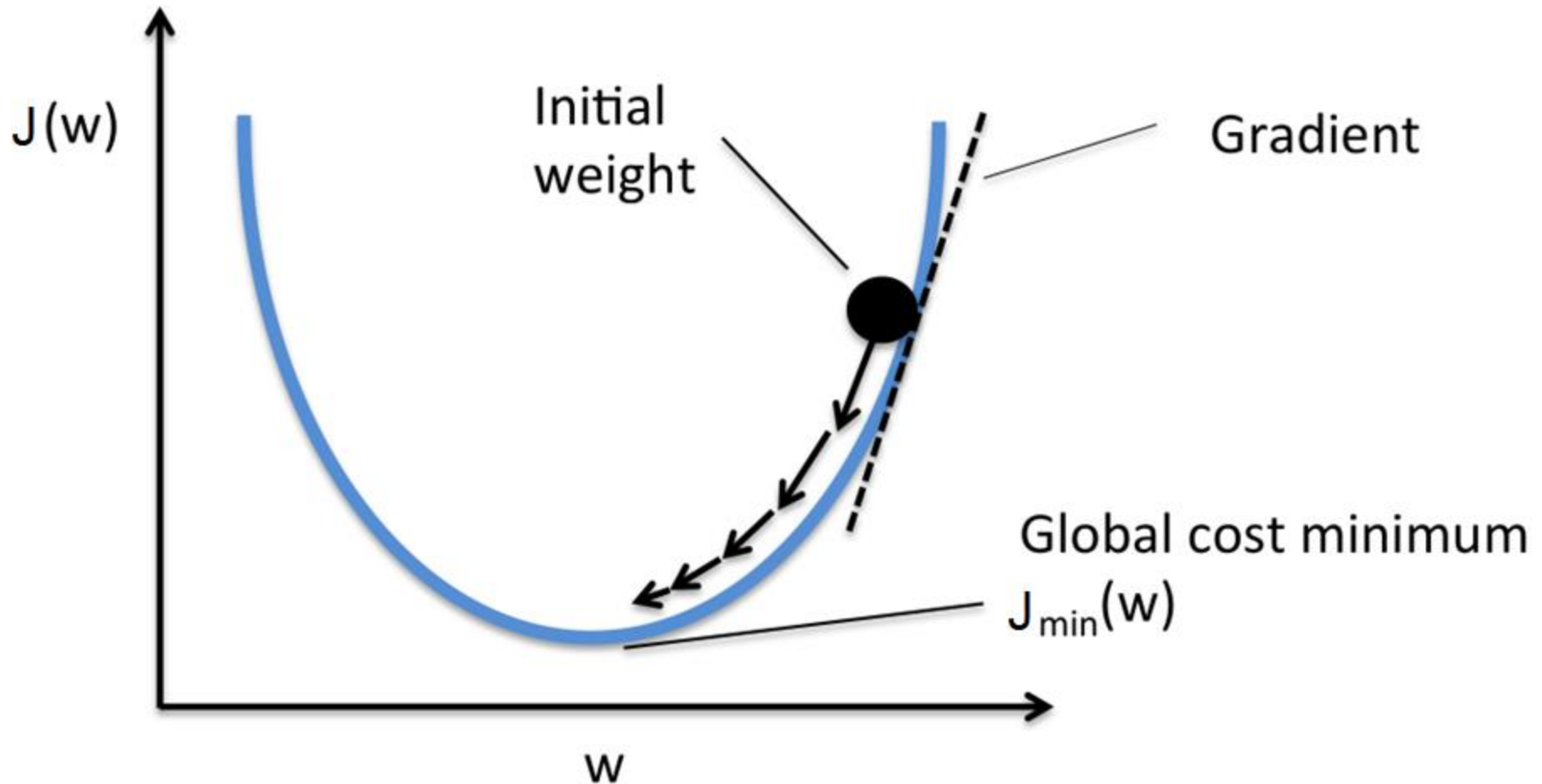
# Tujuan Training

**Meminimalkan Cost Function**

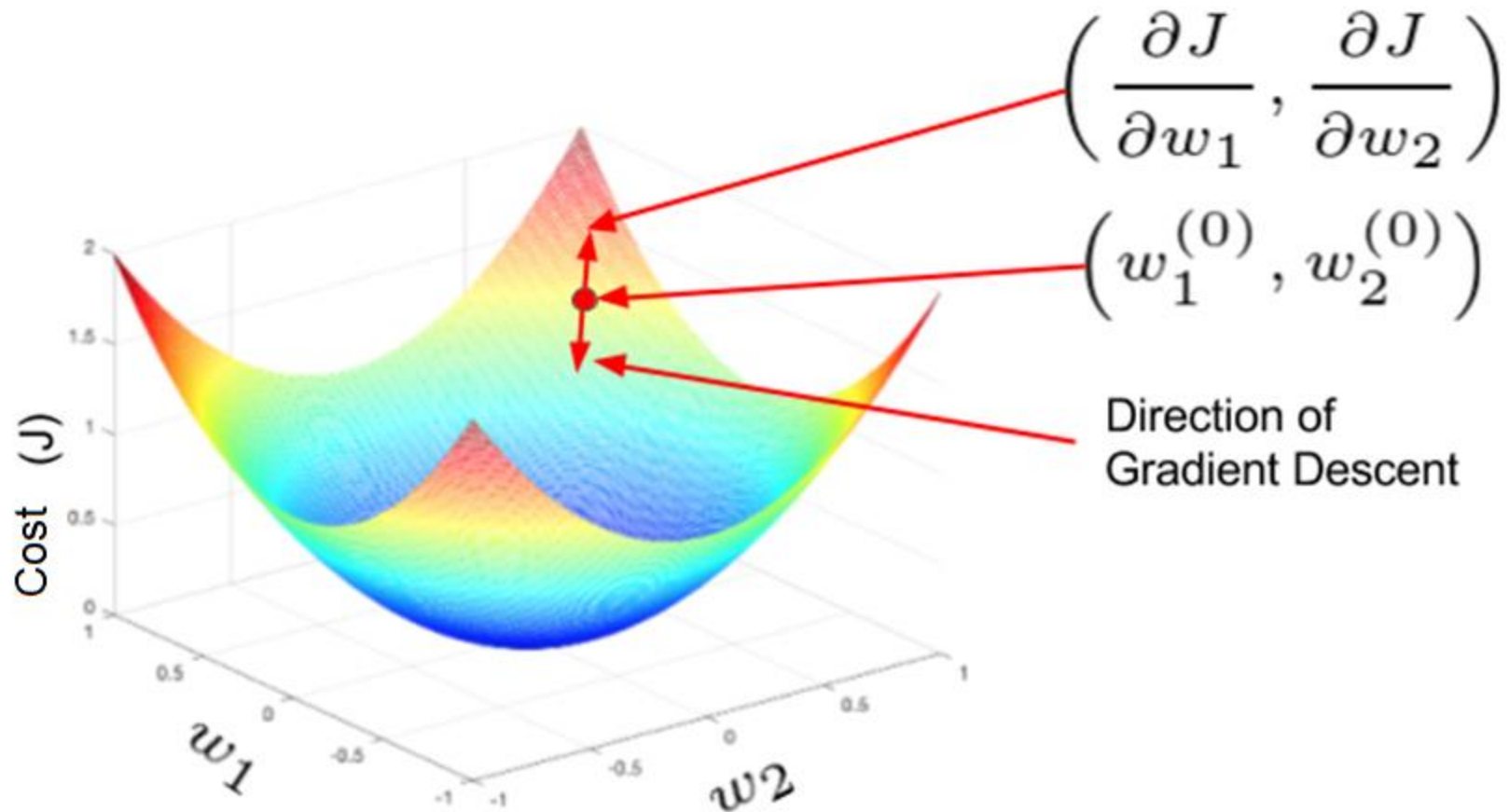


**Dg Merubah Nilai Parameter Model**

# Gradient Descent



# Gradient Descent



# Proses Training

Langkah:

1. Inisialisasi  $w_0$  dan  $w_1$  dengan nilai random, misal:  
 $w_0 = 0.1$  dan  $w_1 = 0.03$
2. Ubah  $w_0$  dan  $w_1$  untuk menurunkan loss/error  $J(w_0, w_1)$
3. Ulangi terus Langkah 2 sampai dengan didapat nilai loss yang paling minimum

Hipotesis:  $h_w(x) = w_0 + w_1x$

Parameter:  $w_0, w_1$

Cost Function:

$$J(w_0, w_1) = \frac{1}{2m} \sum_{i=1}^m (h_w(x^{(i)}) - y^{(i)})^2$$

Goal:  $\underset{w_0, w_1}{\text{minimize}} J(w_0, w_1)$

# Gradient Descent: Mengubah $w_0$ dan $w_1$

*repeat until convergen {*

$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(w)$$

**Update Rule**

Untuk  $j = \{0, \dots, n\}$

*}*

$$temp_0 = w_0 - \alpha \frac{\partial}{\partial w_0} J(w_0, w_1)$$

$$temp_1 = w_1 - \alpha \frac{\partial}{\partial w_1} J(w_0, w_1)$$

$$w_0 = temp_0$$

$$w_1 = temp_1$$

**Keterangan:**

$\alpha$  : Learning Rate

$\frac{\partial}{\partial w_0} J(w_0, w_1)$  : Partial Derivative

# Gradient Descent: Mengubah $w_0$ dan $w_1$

*repeat until convergen {*

*...*

$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(w)$$

Untuk  $j = \{0, \dots, n\}$

*...*

*}*

$$temp_0 = w_0 - \alpha \frac{\partial}{\partial w_0} J(w_0, w_1)$$

$$w_0 = temp_0$$

$$temp_1 = w_1 - \alpha \frac{\partial}{\partial w_1} J(w_0, w_1)$$

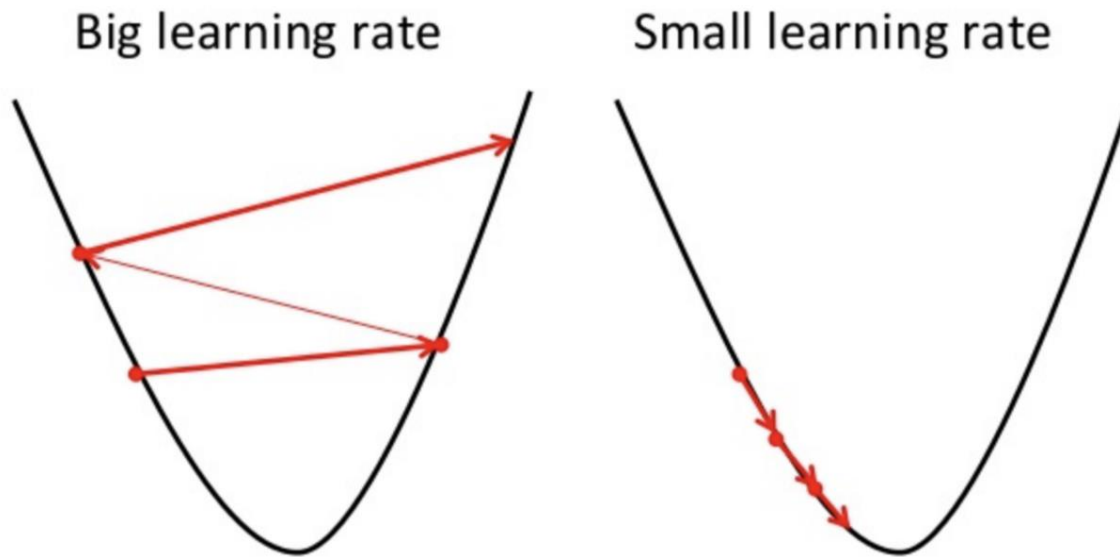
$$w_1 = temp_1$$

**Keterangan:**

$\alpha$  : Learning Rate

$\frac{\partial}{\partial w_0} J(w_0, w_1)$  : Partial Derivative

# Pengaruh Learning Rate



Learning rate menentukan **seberapa besar**(cepat) **perubahan pada bobot**, dkl. **step-size** dari proses update.



# Studi Kasus Sederhana

Hipotesis:

$$h_w(x) = w_0 + w_1x$$

Cost Function:

$$J(w_0, w_1) = \frac{1}{2m} \sum_{i=1}^m (h_w(x^{(i)}) - y^{(i)})^2$$

x	y
2	8
5	17
4	13
7	24
8	26
.	.
.	.
.	.

Update Rule:

$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(w_0, w_1)$$

$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} \frac{1}{2m} \sum_{i=1}^m (h_w(x^{(i)}) - y^{(i)})^2$$

**Ket:**

**i : Index data**

**j : Index parameter**

# Partial Derivatif LMS

$$\frac{\partial J(w)}{\partial w_j} = \frac{\partial}{\partial w_j} \frac{1}{2m} \sum_{i=1}^m (h_w(x^{(i)}) - y^{(i)})^2$$

$$= \frac{1}{2m} \sum_{i=1}^m \frac{\partial}{\partial w_j} (h_w(x^{(i)}) - y^{(i)})^2$$

$$= \frac{1}{2m} \sum_{i=1}^m \frac{\partial}{\partial w_j} (\mathbf{h}_w(\mathbf{x}^{(i)}) - \mathbf{y}^{(i)}) (\mathbf{h}_w(\mathbf{x}^{(i)}) - \mathbf{y}^{(i)})$$

$$= \frac{1}{2m} \sum_{i=1}^m 2(h_w(x^{(i)}) - y^{(i)}) \frac{\partial}{\partial w_j} (h_w(x^{(i)}) - y^{(i)})$$

$$= \frac{1}{2m} \sum_{i=1}^m 2(h_w(x^{(i)}) - y^{(i)}) (x_j^{(i)}) = \frac{1}{m} \sum_{i=1}^m (h_w(x^{(i)}) - y^{(i)}) (x_j^{(i)})$$

jika  $f(x) = \mathbf{u}(x) \mathbf{v}(x)$

maka  $f'(x) = \mathbf{u}'(x) \mathbf{v}(x) + \mathbf{u}(x) \mathbf{v}'(x)$

$$h_w(x) = w_0 + w_1 x$$

# Studi Kasus Sederhana

Hipotesis:

$$h_w(x) = w_0 + w_1 x$$

Cost Function:

$$J(w_0, w_1) = \frac{1}{2m} \sum_{i=1}^m (h_w(x^{(i)}) - y^{(i)})^2$$

x	y
2	8
5	17
4	13
7	24
8	26
.	.

Update Rule:

$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(w_0, w_1) = w_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_w(x^{(i)}) - y^{(i)}) (x_j^{(i)})$$

$$w_0 = w_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_w(x^{(i)}) - y^{(i)}) (x_0^{(i)})$$

$$w_1 = w_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_w(x^{(i)}) - y^{(i)}) (x_1^{(i)})$$

**Ket:**

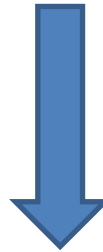
**i : Index data**

**j : Index parameter**

# Gradient Descent

# Tujuan Training

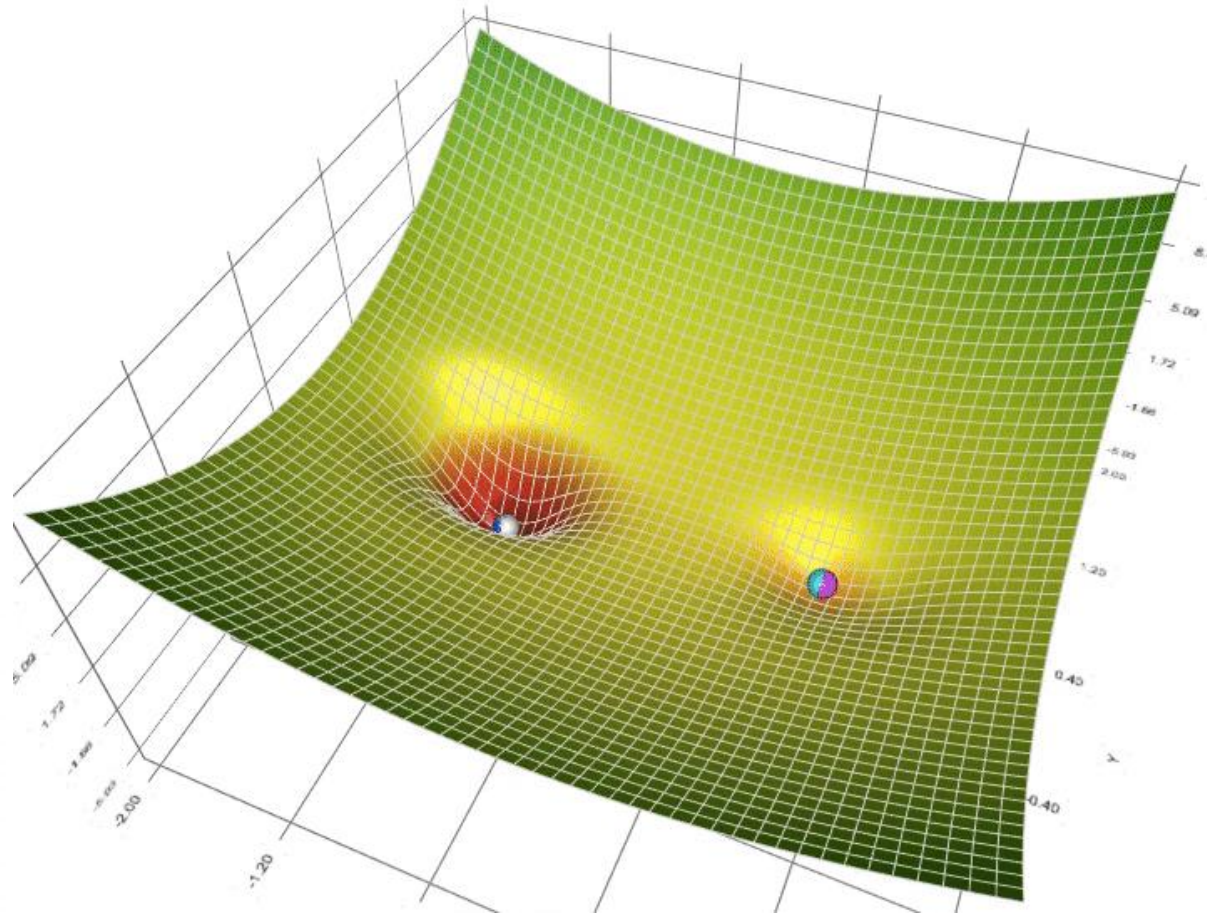
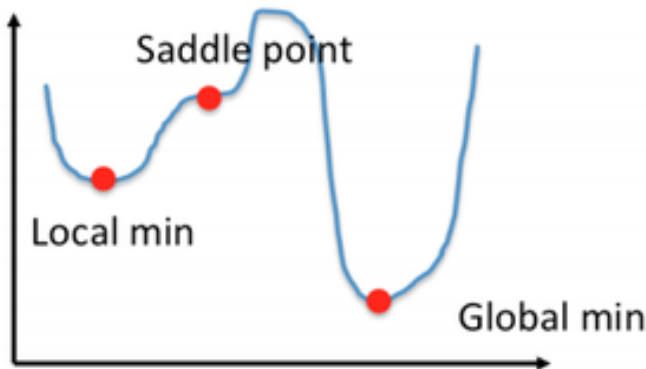
**Meminimalkan Nilai Error**



**Dg Merubah Nilai Bobot**

# Requirement Algoritma Learning

- Menemukan global minima dengan cepat dan reliabel
- Tidak terjebak di local minima, saddle point, atau bahkan pada plateu

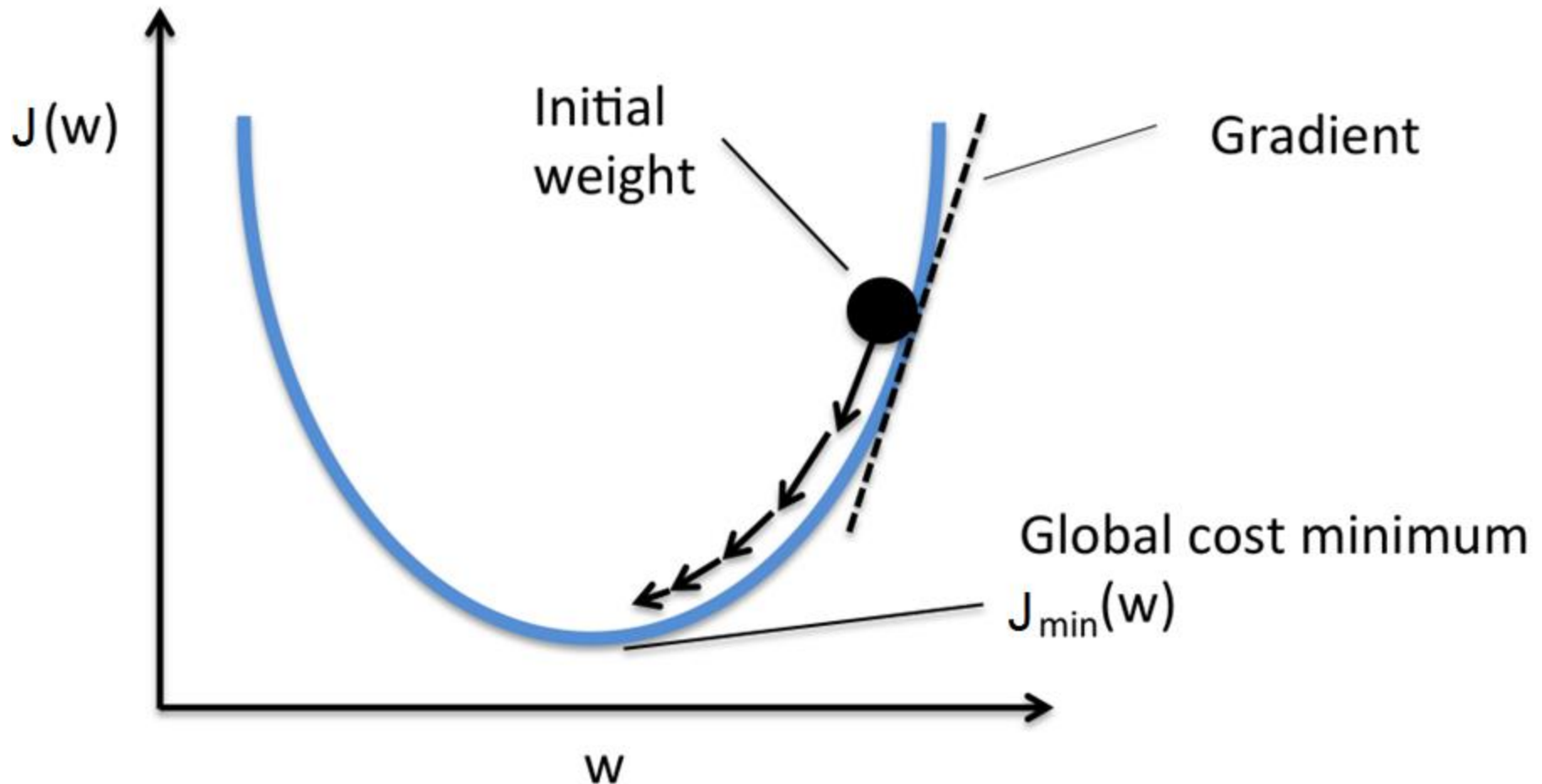


# Gradient Descent

- Salah satu algoritma optimasi yg dapat digunakan utk meminimalkan fungsi error berdasarkan bobot pada JST
- **Learning rate**: factor pengali yang menentukan seberapa besar nilai update untuk suatu bobot (i.e.: 0.1)
- Satu **epoch**: Satu Langkah lengkap yang memproses seluruh data training
- Gradien menunjukkan tingkat kemiringan suatu fungsi pada suatu titik tertentu (x,y), dapat dicari dengan menggunakan turunan pertama fungsi tersebut terhadap x.
- Contoh fungsi error sederhana (linear):  
$$\text{Error} = Y(\text{Predicted}) - T(\text{Target/Actual})$$

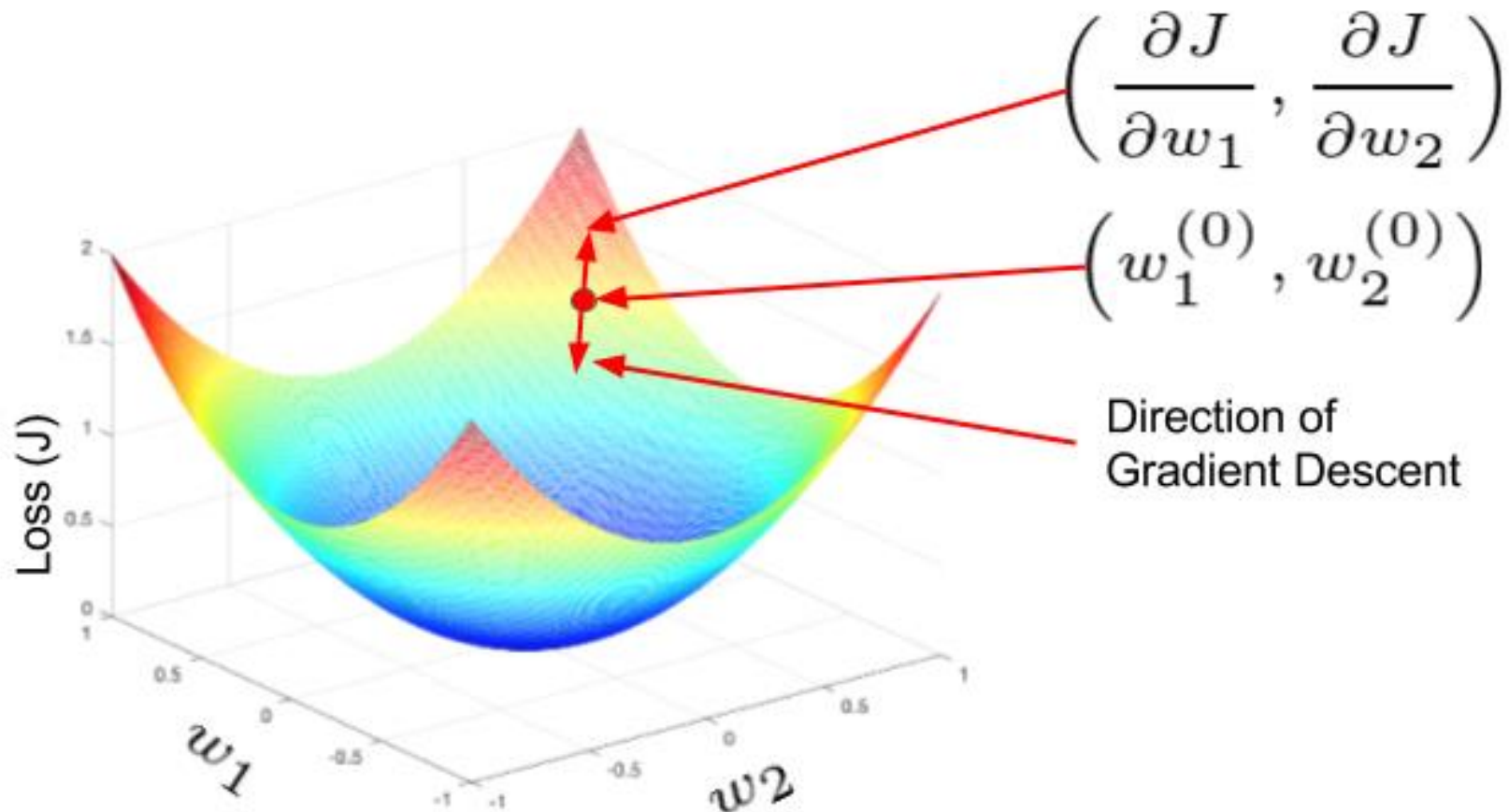
# Minimasi Fungsi Error dg Merubah Bobot

Merubah nilai initial  $w \rightarrow$  Menuju titik  $J$  minimum  $\rightarrow$  Dipandu oleh gradient





# Minimasi: w lebih dari satu var/fitur



# Backprop dan Updating Rule

Meminimalkan Nilai Error



Dg Merubah Nilai Bobot

## Gradien Negatif:

( pada posisi  $w_{i1}^k$  )

Jika  $w$  ditambah, maka Loss function menurun.

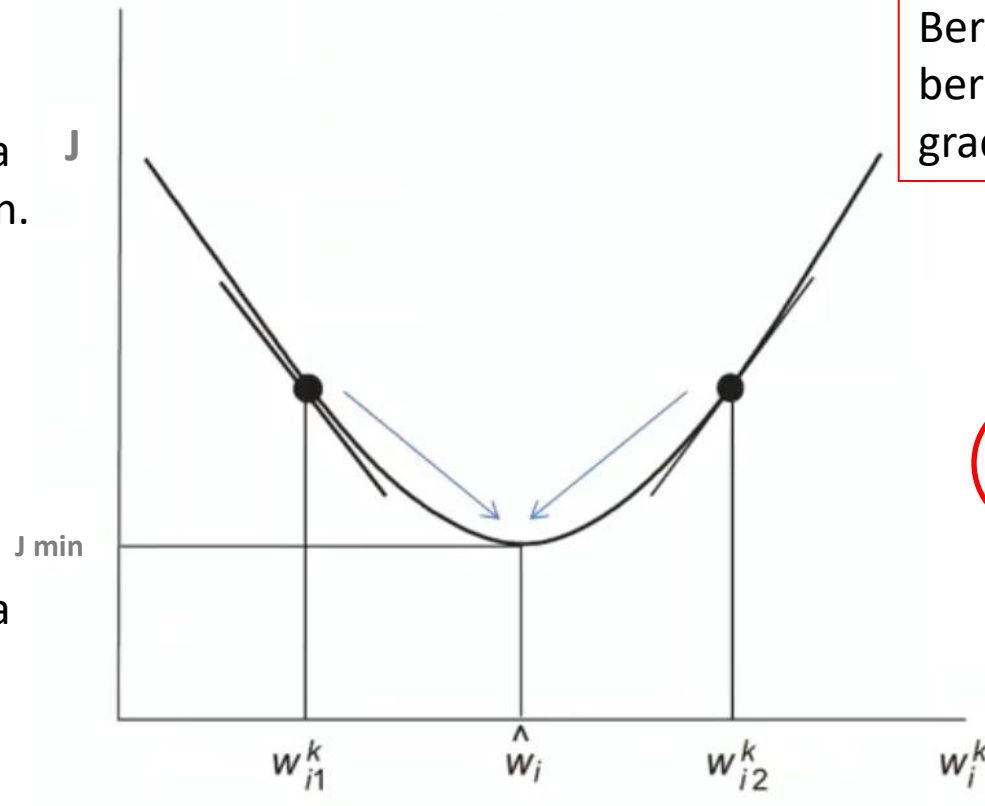
$-(-m) = + \rightarrow w$  harus dinaikan/ditambah (bergerak ke kanan)

## Gradien Positif:

( pada posisi  $w_{i2}^k$  )

Jika  $w$  ditambah, maka Loss function naik.

$-(+m) = - \rightarrow w$  harus diturunkan/dikurangi (bergerak ke kiri)



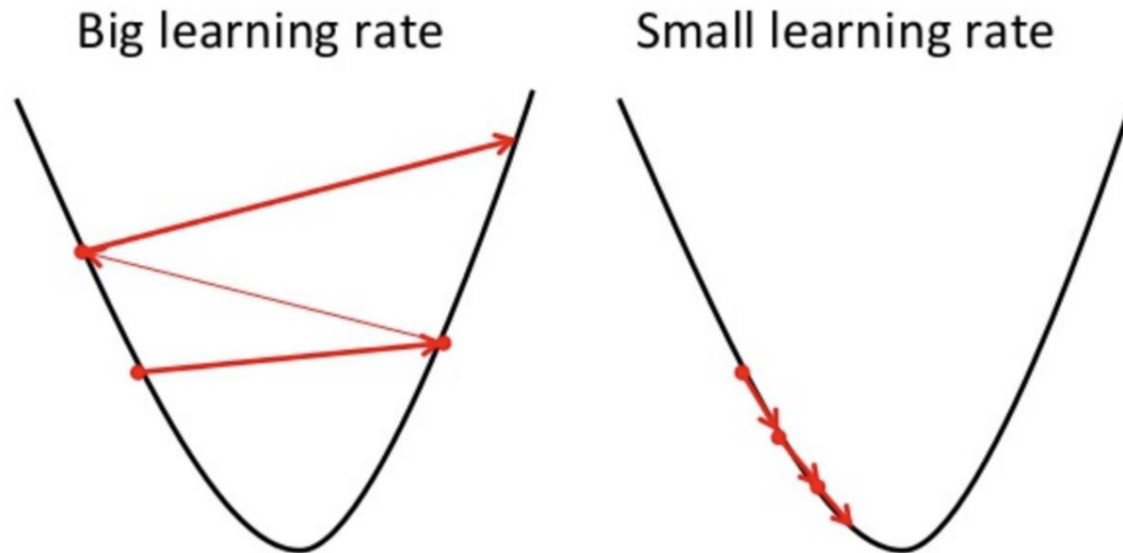
Bergerak ke arah yg berlawanan dengan gradien



$$-\frac{\partial J}{\partial w_i^k}$$

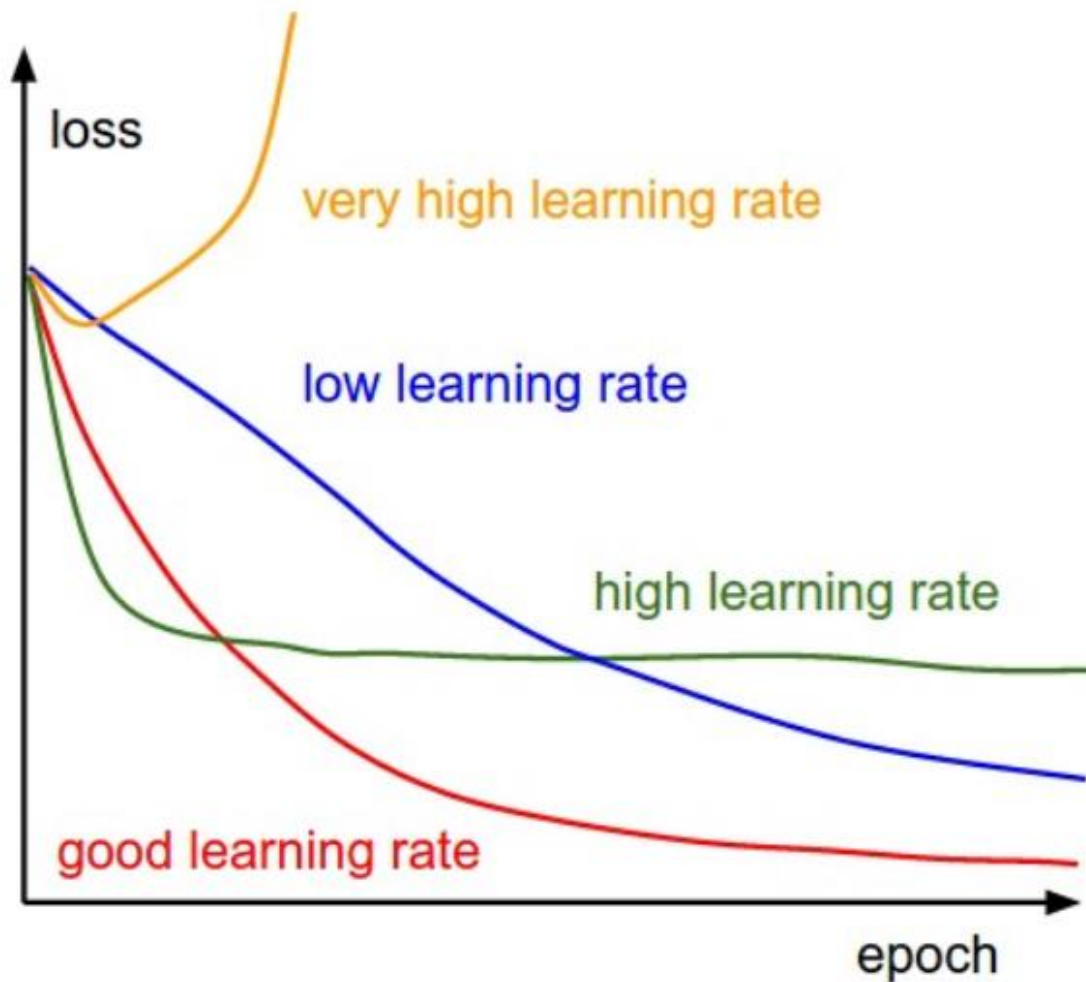
# Weight Update Rule

bobot baru		bobot lama	negasi dari gradien	learning rate	gradient
$\mathbf{w}$	=	$\mathbf{w}$	—	$\alpha$	$\frac{\partial J}{\partial \mathbf{w}}$



Learning rate menentukan **seberapa besar**(cepat) **perubahan pada bobot**, dkl. **step-size** dari proses update.

# Pengaruh Learning Rate



Learning rate yg banyak digunakan:  $0.001$ ,  $0.003$ ,  $0.01$ ,  $0.03$ ,  $0.1$ ,  $0.3$

# (Vanilla) Gradient Descent: Algoritma

1. Inisialisasi bobot secara random
2. Hitung output (o) dari model
3. Hitung gradient:  $\frac{\partial J}{\partial w}$
4. Update bobot menggunakan Weight Update Rule
5. Cek apakah kondisi berhenti terpenuhi? Jika ya, selesai. Jika tidak, ulangi Langkah 2.

Kondisi berhenti:

- Slope ditemukan (missal: error  $\leq 0.001$ )
- Maksimum epoch tercapai
- Konvergen  $\rightarrow$  Error tidak berubah (flat)
- Dll.

# Aturan Pembelajaran

Gradient  $\nabla J[\vec{w}] \equiv \frac{\partial J}{\partial w} \equiv \left[ \frac{\partial J}{\partial w_0}, \frac{\partial J}{\partial w_1}, \dots, \frac{\partial J}{\partial w_n} \right]$

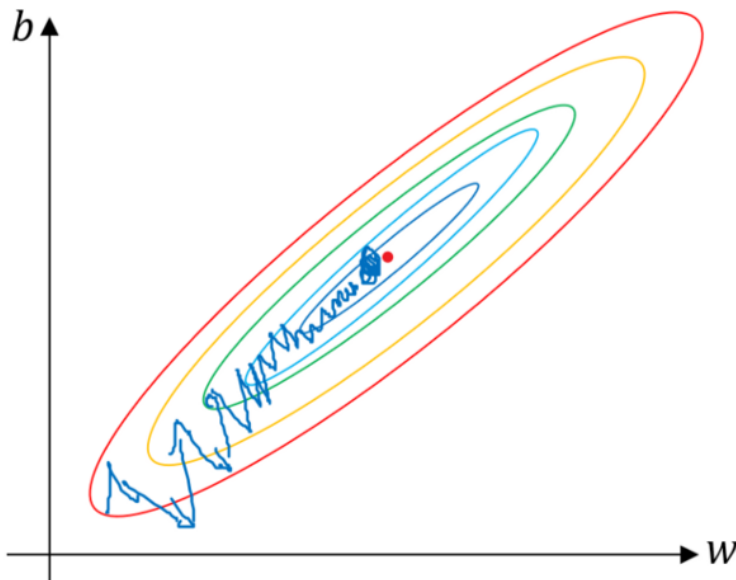
Training rule:  $\Delta w = -\alpha \nabla J[\vec{w}]$

*maka*  $\Delta w_i = -\alpha \frac{\partial J}{\partial w_i}$

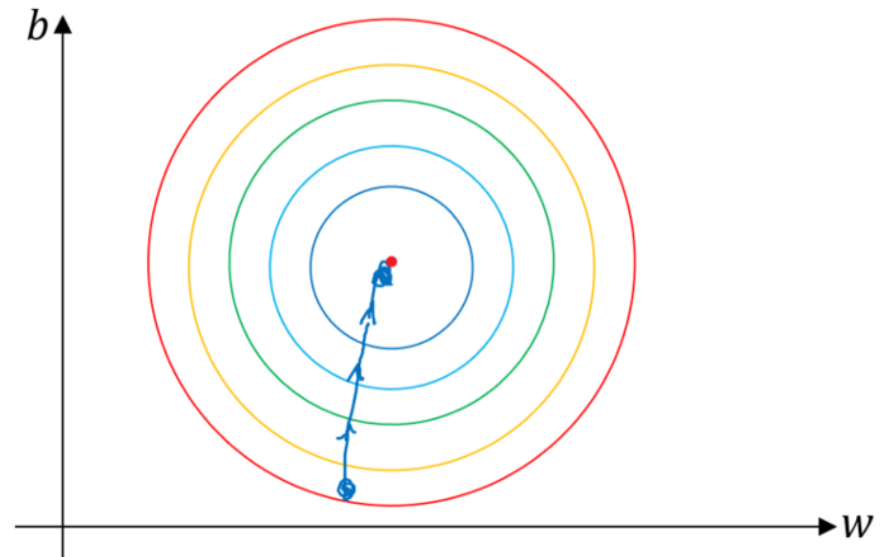
# Normalisasi Data

- Data yang digunakan harus dinormalisasi
- Pada GD ini cocok digunakan scaling data menjadi range antara 0 s.d. 1

Unnormalized



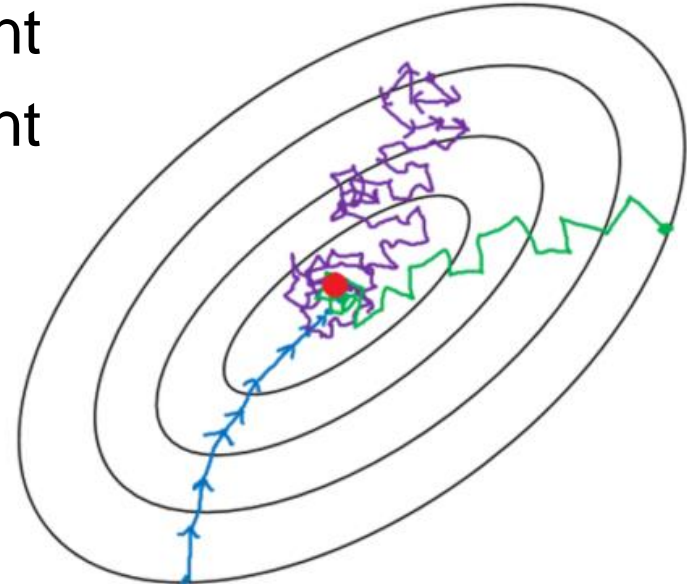
Normalized



# Strategi Update Bobot

- Variasi: Berbeda pada jumlah data yg digunakan untuk menghitung gradient pada setiap epoch-nya.
- Trade-off: Akurasi vs Waktu Eksekusi
- Variasi Strategi:
  - Batch Gradient Descent
  - Mini-batch Gradient Descent
  - Stochastic Gradient Descent

— Batch gradient descent  
— Mini-batch gradient Descent  
— Stochastic gradient descent





# Batch Gradient Descent

- Untuk melakukan update bobot, harus menggunakan semua data dalam proses perhitungan gradient-nya

```
for i in range(num_epochs):  
    grad = compute_gradient(data, bobot)  
    bobot = bobot - learning_rate * grad
```

- Dalam setiap epoch, hanya sekali melakukan update bobot
- **Kelebihan:**
  - Stabil, convergence
  - Grafik fungsi loss lebih smooth
- **Kekurangan:**
  - Butuh memory yang besar
  - Waktu yang dibutuhkan lama, apalagi datanya banyak

# Mini-Batch Gradient Descent

- Untuk melakukan update bobot, hanya menggunakan sebagian data (disebut **batch**) dalam proses perhitungan gradient-nya

```
for i in range(num_epochs):  
    np.random.shuffle(data)  
    for batch in radom_minibatches(data, batch_size=32):  
        grad = compute_gradient(batch, bobot)  
        bobot = bobot - learning_rate * grad
```

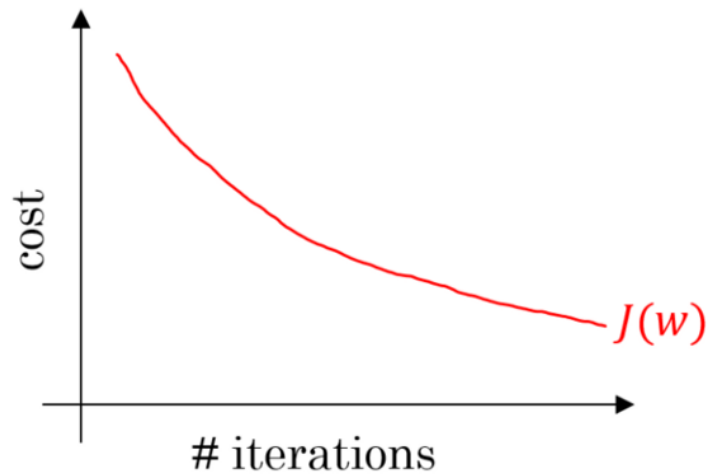
- Setiap epoch, dilakukan (**jml / jml-per-batch**) kali update bobot
- Misal jml data = 10.000, jml per batch = 32, maka dalam setiap epoch dilakukan  $10.000/32 = 313$  kali update
- Jika jml data dibagi jml data per batch tidak bulat, maka sisa data pada partisi terakhir dianggap satu batch
- Jika jml data per batch sama dengan jml data, maka sama dg batch-GD

# Mini-Batch Gradient Descent

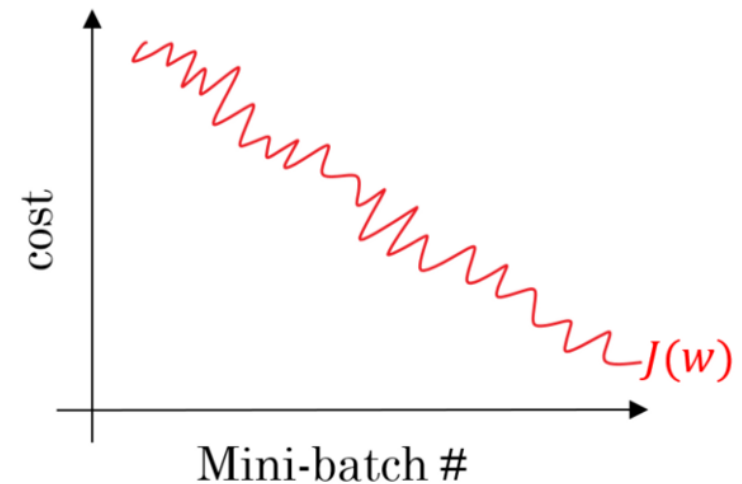
- **Kelebihan:**
  - Lebih cepat waktu belajarnya dibanding dg batch GD
  - Ada tambahan noise, sehingga meningkatkan generalization
  - Lebih ramah memory, karena dalam satu proses hanya mengolah sejumlah kecil data
- **Kekurangan:**
  - Sulit convergence, karena ada tambahan noise
  - Grafik Loss lebih fluktuatif

# Batch vs Mini-Batch GD

Batch gradient descent



Mini-batch gradient descent



# Berapa Batch Size yang Tepat?

Jml data per batch biasanya merupakan hasil perpangkatan dengan basis 2: 32, 64, 128, 256, 512, dll, karena GPU performanya lebih optimal pada nilai-nilai tsb.

← Tweet



Yann LeCun  
@ylecun

...

Training with large minibatches is bad for your health.  
More importantly, it's bad for your test error.  
Friends dont let friends use minibatches larger than 32.



arxiv.org

Revisiting Small Batch Training for Deep Neural Networks  
Modern deep neural network training is typically based on  
mini-batch stochastic gradient optimization. While the use ...

4:00 AM · Apr 27, 2018 · Facebook

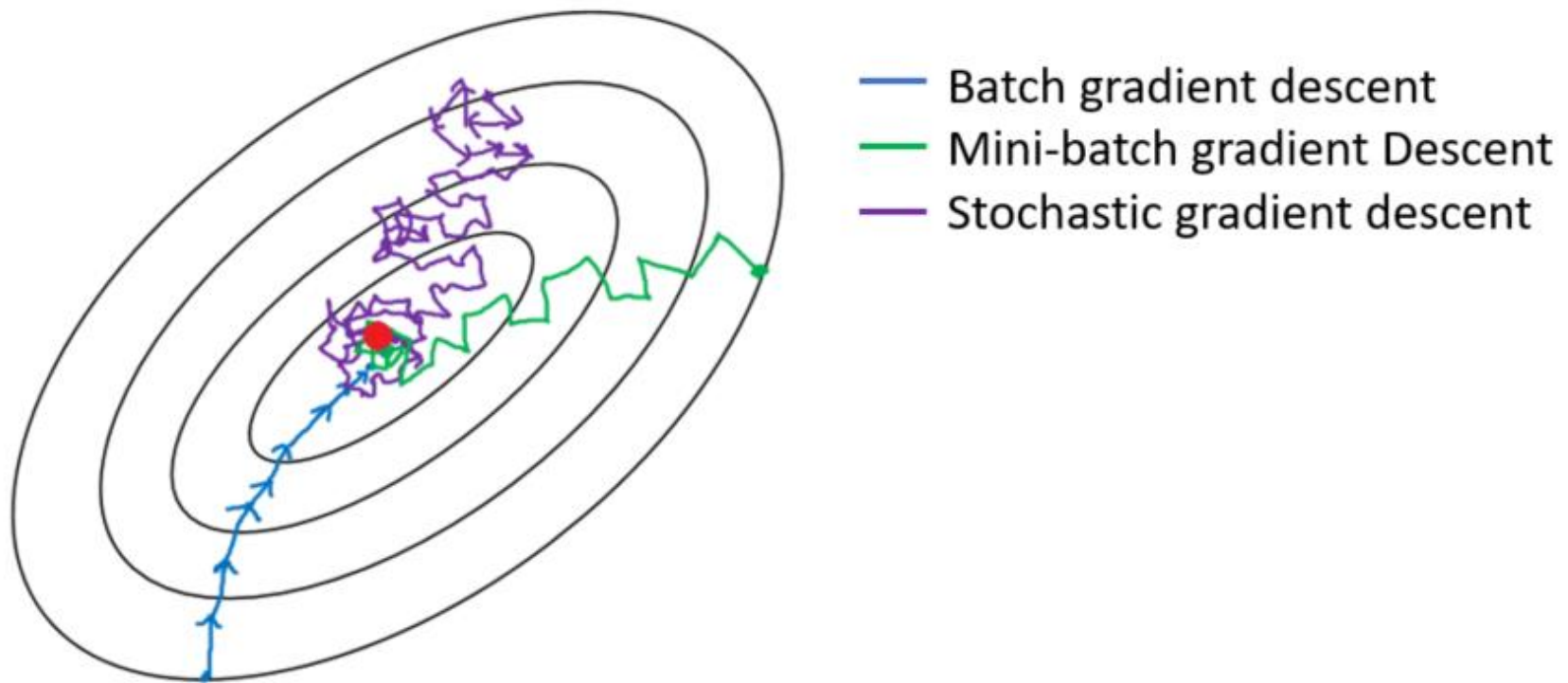
# Stochastic Gradient Descent

- Update bobot dilakukan untuk setiap data, artinya gradient dihitung untuk setiap data sample

```
for i in range(num_epochs):  
    np.random.shuffle(data)  
    for example in data:  
        grad = compute_gradient(example, bobot)  
        bobot = bobot - learning_rate * grad
```

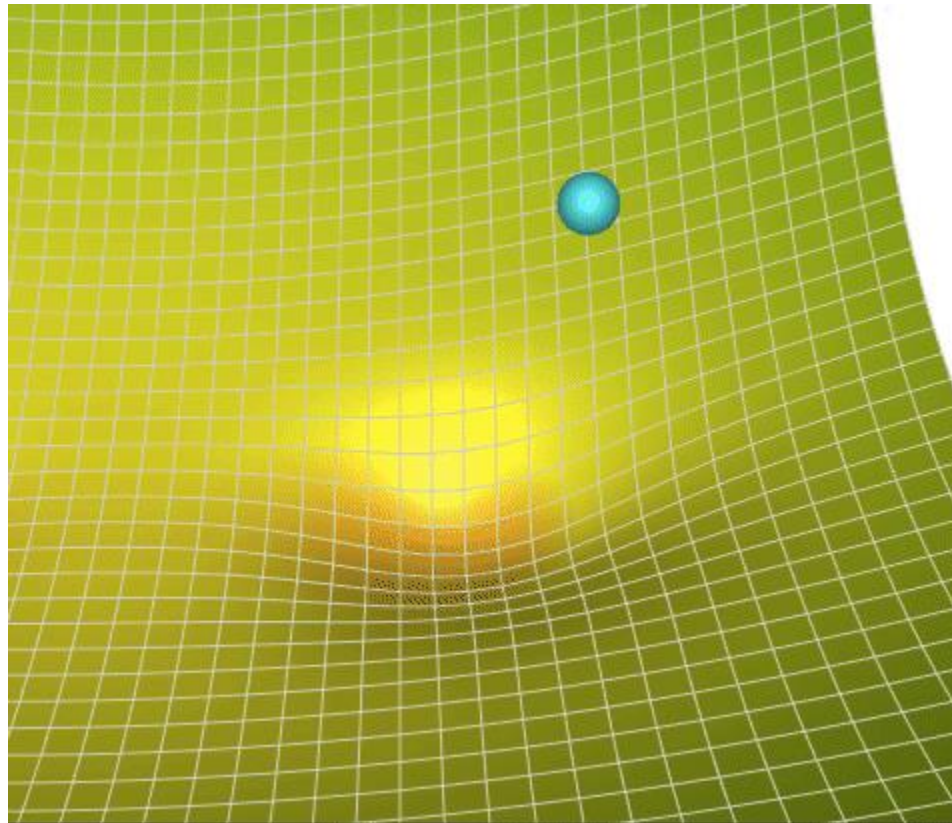
- Setiap epoch, dilakukan update bobot sebanyak jumlah data
- Hemat memory, komputasi cepat, efektif utk jml data yang besar
- Namun, noisenya sangat tinggi

# Perbandingan Batch, Mini-Batch, dan Stochastic



# Gradient Descent Improvements

**Vanilla Gradient Descent lambat dan sering terjebak pada local minima**





# Gradient Descent Improvements

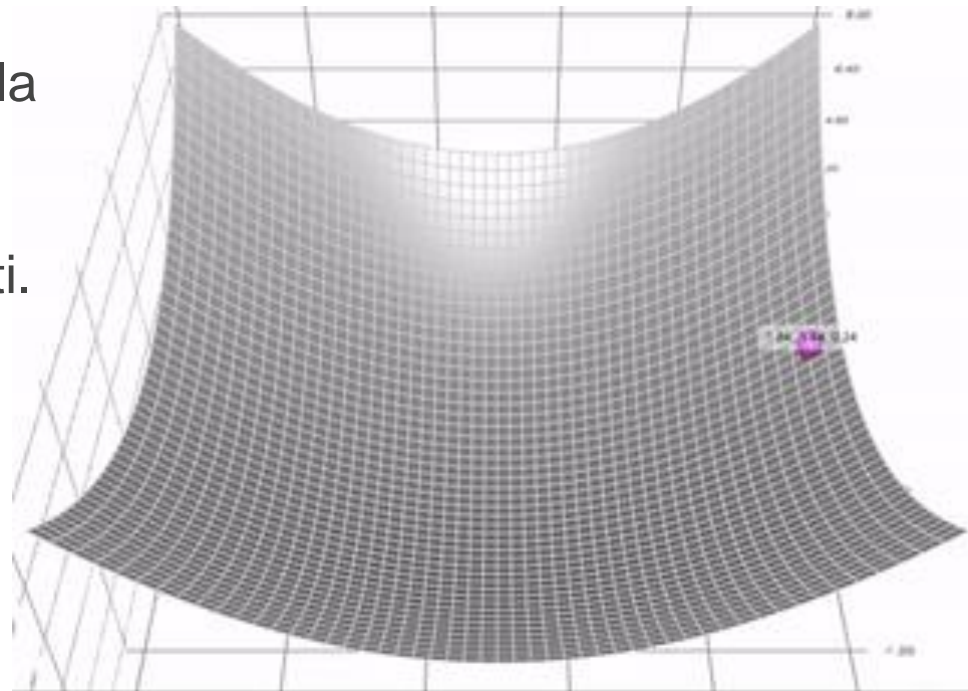
- Beberapa upaya peningkatan GD:
  - Momentum Based Gradient (Momentum)
  - Adaptive Gradient (AdaGrad)
  - Root Mean Square Propagation (RMSProp)
  - Adaptive Moment Estimation (Adam)

# Momentum Based Gradient (Momentum)

- Menerapkan konsep momentum dalam ilmu Fisika
- Momentum adalah besaran yang dihasilkan dari jumlah perkalian skala massa benda dan vektor kecepatan gerak benda

$$\mathbf{p} = m \mathbf{v}$$

- Bayangkan sebuah bola digelindingkan pada sebuah mangkuk dengan zero gesekan, dengan adanya momentum, bola tersebut alih-alih berhenti pada cekungan, malah akan terus bergerak bolak-balik tanpa henti.
- Konsep ini dapat diterapkan pada vanilla GD

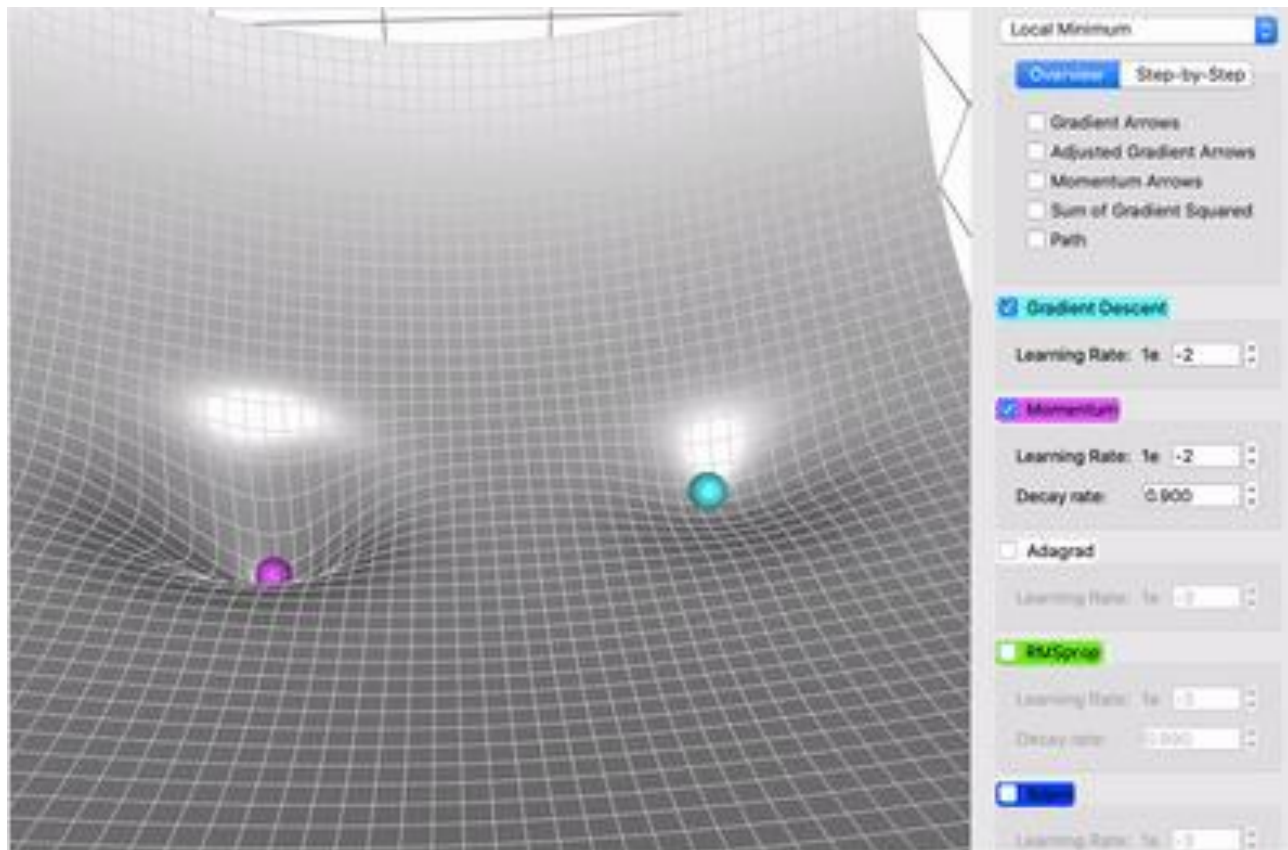


# Momentum Based Gradient (Momentum)

- Vanilla GD:  
 $\text{delta} = - (\text{learning\_rate} * \text{gradient})$   
 $\text{bobot} = \text{bobot} + \text{delta}$
- Momentum  
 $\text{delta} = - (\text{learning\_rate} * \text{gradient}) + \text{previous\_delta} * \text{decay\_rate}$   
 $\text{bobot} = \text{bobot} + \text{delta}$
- Faktorisasi:  
 $\text{sum\_of\_grad} = \text{gradient} + \text{previous\_sum\_of\_grad} * \text{decay\_rate}$   
 $\text{delta} = -\text{learning\_rate} * \text{sum\_of\_grad}$   
 $\text{bobot} = \text{bobot} + \text{delta}$
- $\text{decay\_rate} = 0 \rightarrow$  Vanilla Gradient Descent
- $\text{decay\_rate} = 1 \rightarrow$  Tak akan pernah berhenti
- $\text{decay\_rate}$  ideal  $\rightarrow 0.8$  s.d.  $0.9 \rightarrow$  perlahan akan berhenti

# Momentum Based Gradient (Momentum)

- Momentum membuat pergerakan menjadi lebih cepat
- Mampu menghindari/keluar dari jebakan local minima
- Momentum (magenta) vs. Gradient Descent (cyan):

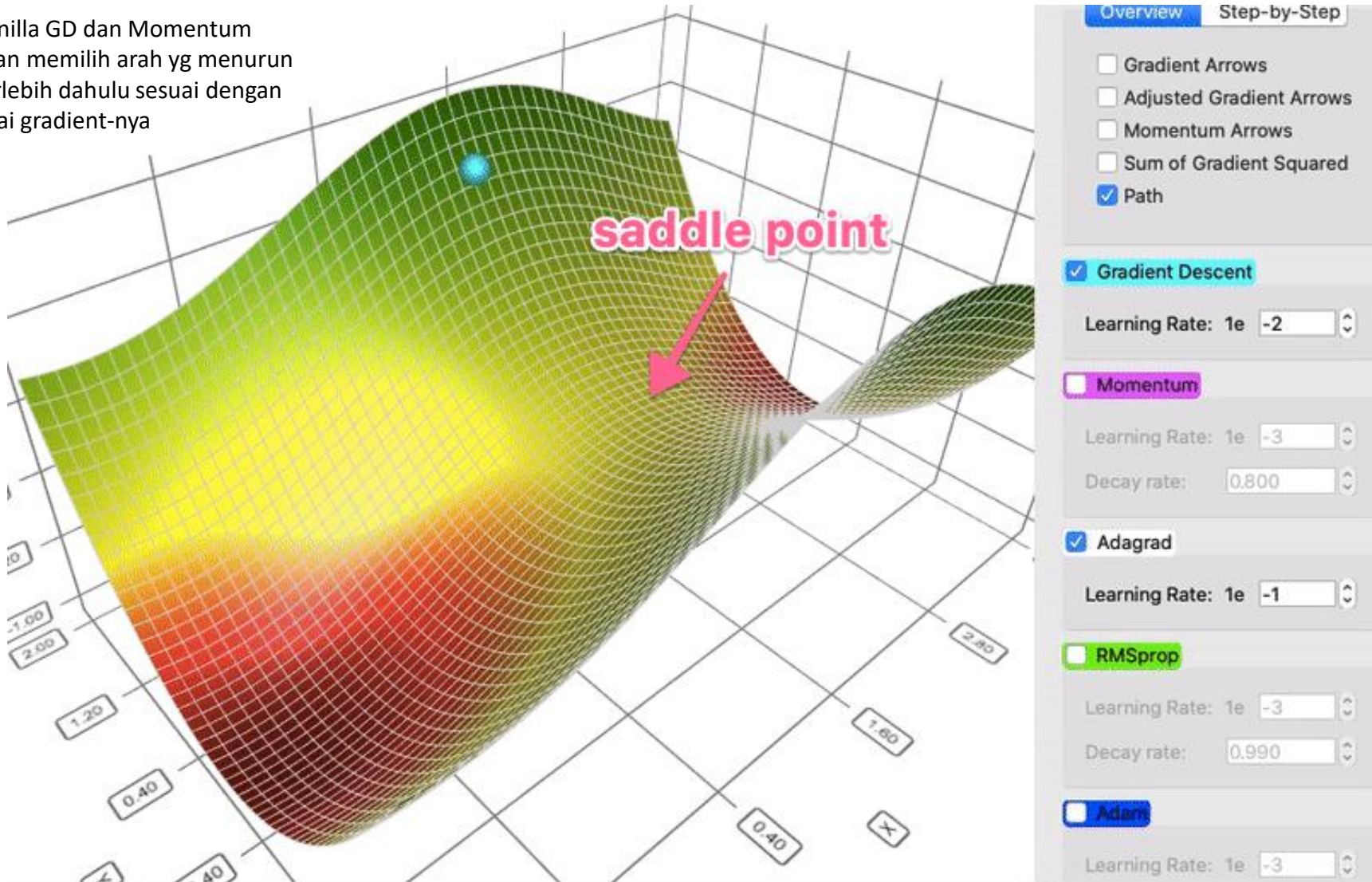


# Adaptive Gradient Descent (AdaGrad)

- Ada kalanya fitur yang digunakan bersifat sparse
- Rataan dari gradient fitur yang sparse cenderung kecil, berakibat pada lamanya proses training, karena pergerakannya lambat.
- AdaGrad mengatasi ini dengan ide: Semakin besar update pada bobot, maka akan semakin kecil update selanjutnya, hal ini memberi kesempatan pada fitur yang sparse untuk melakukan penyesuaian.
- Pada implementasinya, AdaGrad menggunakan akumulasi kuadrat gradient untuk memandu arah pergerakan ini.
- Hal ini membuat AdaGrad (dan variasinya) mampu keluar dari jebakan saddle point.
- AdaGrad:  
$$\begin{aligned} \text{sum\_of\_grad\_squared} &= \text{previous\_sum\_of\_grad\_squared} + \text{gradient}^2 \\ \text{delta} &= -\text{learning\_rate} * \text{gradient} / \sqrt{\text{sum\_of\_grad\_squared}} \\ \text{bobot} &= \text{bobot} + \text{delta} \end{aligned}$$

# AdaGrad (white) vs. gradient descent (cyan)

Vanilla GD dan Momentum akan memilih arah yg menurun terlebih dahulu sesuai dengan nilai gradient-nya

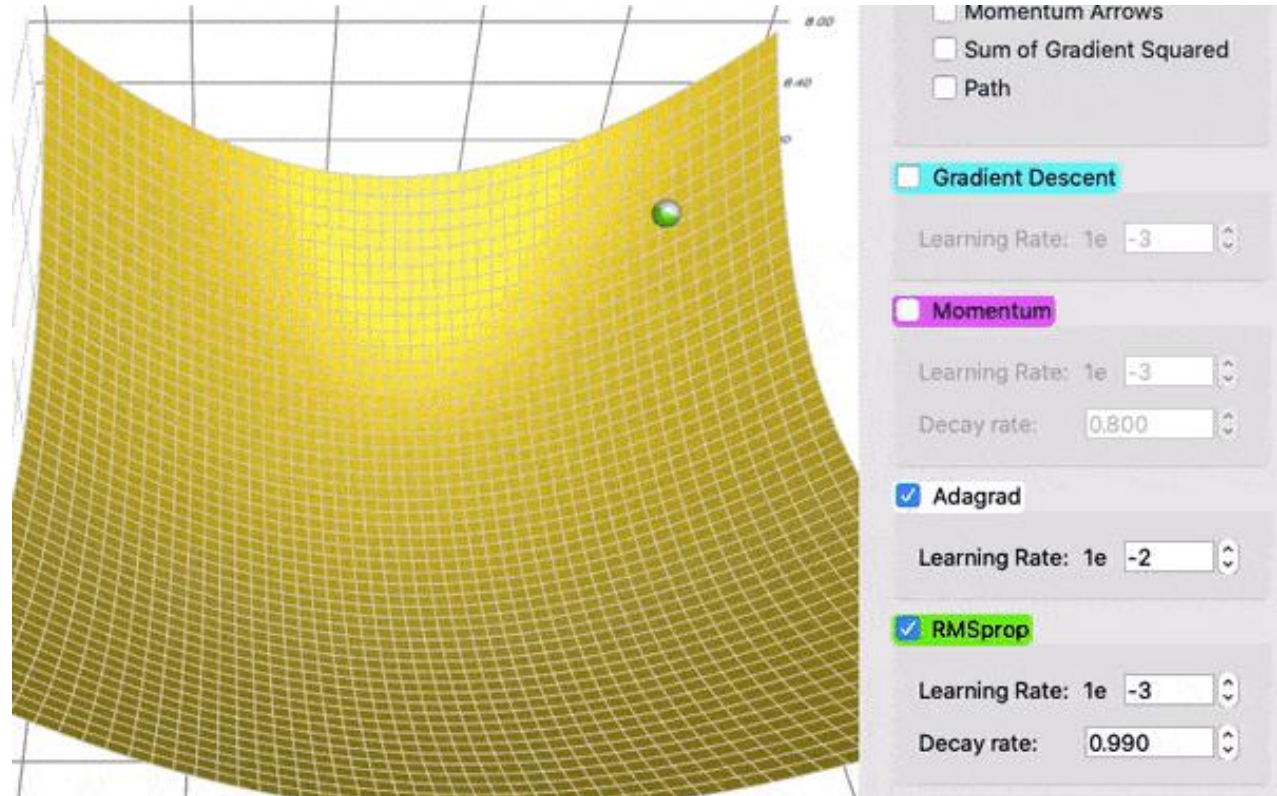




# Root Mean Square Propagation (RMSProp)

- Kekurangan dari AdaGrad adalah lambat, karena `sum_of_gradient_squared` hanya akan bertambah terus menerus tanpa penurunan.
- RMSProp mengatasinya dengan menambahkan `decay_rate`

RMSProp (green)  
vs AdaGrad (white)



- RMSProp:  

$$\text{sum\_of\_grad\_squared} = \text{previous\_sum\_of\_grad\_squared} * \text{decay\_rate} + \text{gradient}^2 * (1 - \text{decay\_rate})$$

$$\text{delta} = -\text{learning\_rate} * \text{gradient} / \sqrt{\text{sum\_of\_grad\_squared}}$$

$$\text{bobot} = \text{bobot} + \text{delta}$$

# Adaptive Moment Estimation (Adam)

- Mengabungkan keunggulan RMSProp dan Momentum

- Adam:

$\text{sum\_of\_grad} = \text{previous\_sum\_of\_grad} * \text{beta1} + \text{gradient} * (1 - \text{beta1})$

[Momentum]

$\text{sum\_of\_grad\_squared} = \text{previous\_sum\_of\_grad\_squared} * \text{beta2} + \text{gradient}^2 * (1 - \text{beta2})$

[RMSProp]

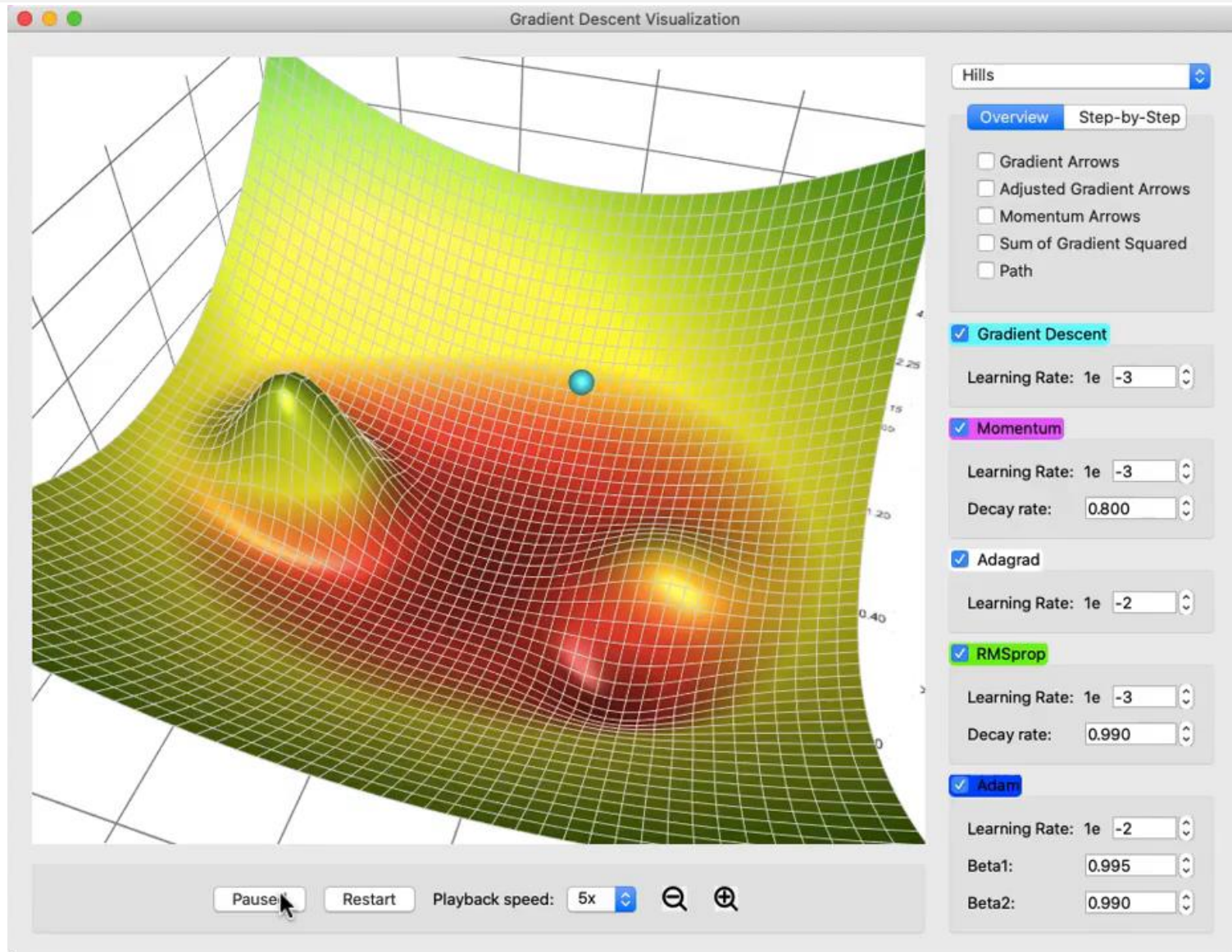
$\text{delta} = -\text{learning\_rate} * \text{sum\_of\_gradient} / \sqrt{\text{sum\_of\_grad\_squared}}$

$\text{bobot} = \text{bobot} + \text{delta}$

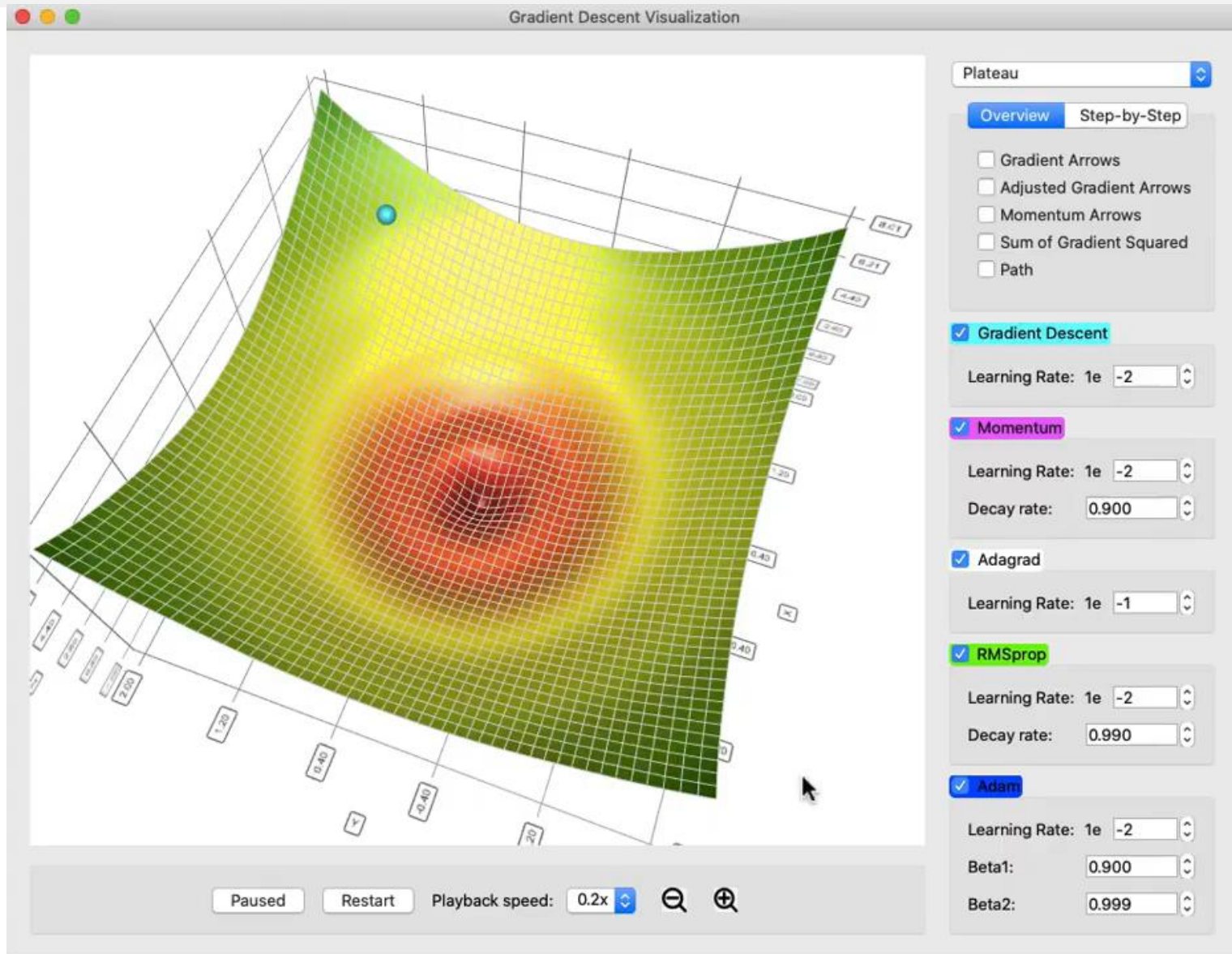
- Beta1 merupakan decay\_rate untuk moment pertama (biasanya: 0.9)
- Beta2 merupakan decay\_rate untuk moment kedua (biasanya: 0.999)
- Adam mewarisi kecepatan dari momentum dan RMSProp sebagai pemandu arah gerakannya



# Viz: GD on Hills



# Viz: GD on Plateau



# Thank You

