# Gradient Descent

Guowei Wei
Department of Mathematics
Michigan State University

*References:*
*Duc D. Nguyen's lecture notes*
*Wikipedia*

**Artificial Intelligence (AI) Stats News (**Sep 10, 2019**):**
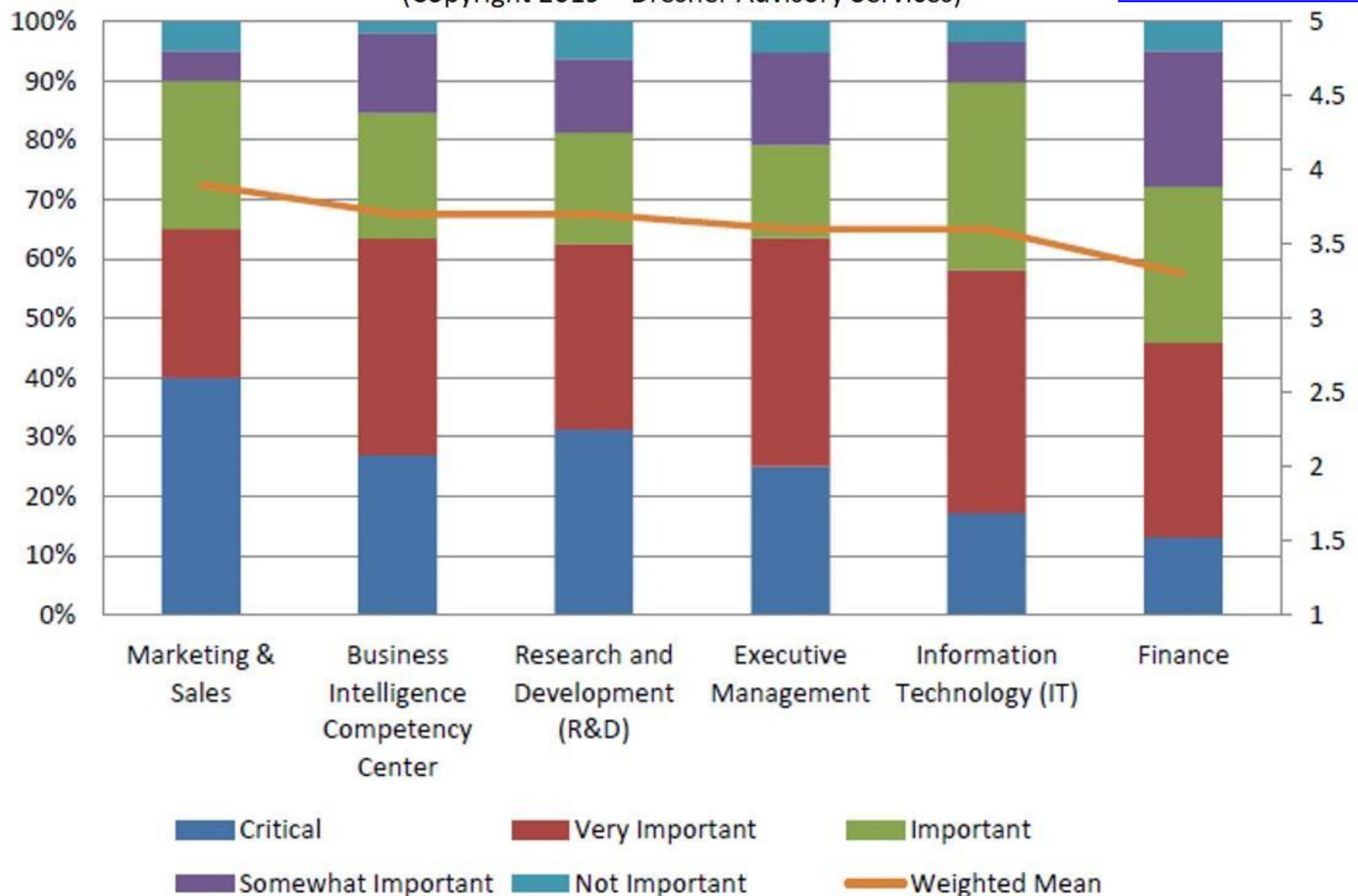**120 Million Workers Need To Be Retrained Because Of AI in the next three years**



**Gil Press**

# Importance of Data Science and Machine Learning by Function

(Copyright 2019 – Dresner Advisory Services)
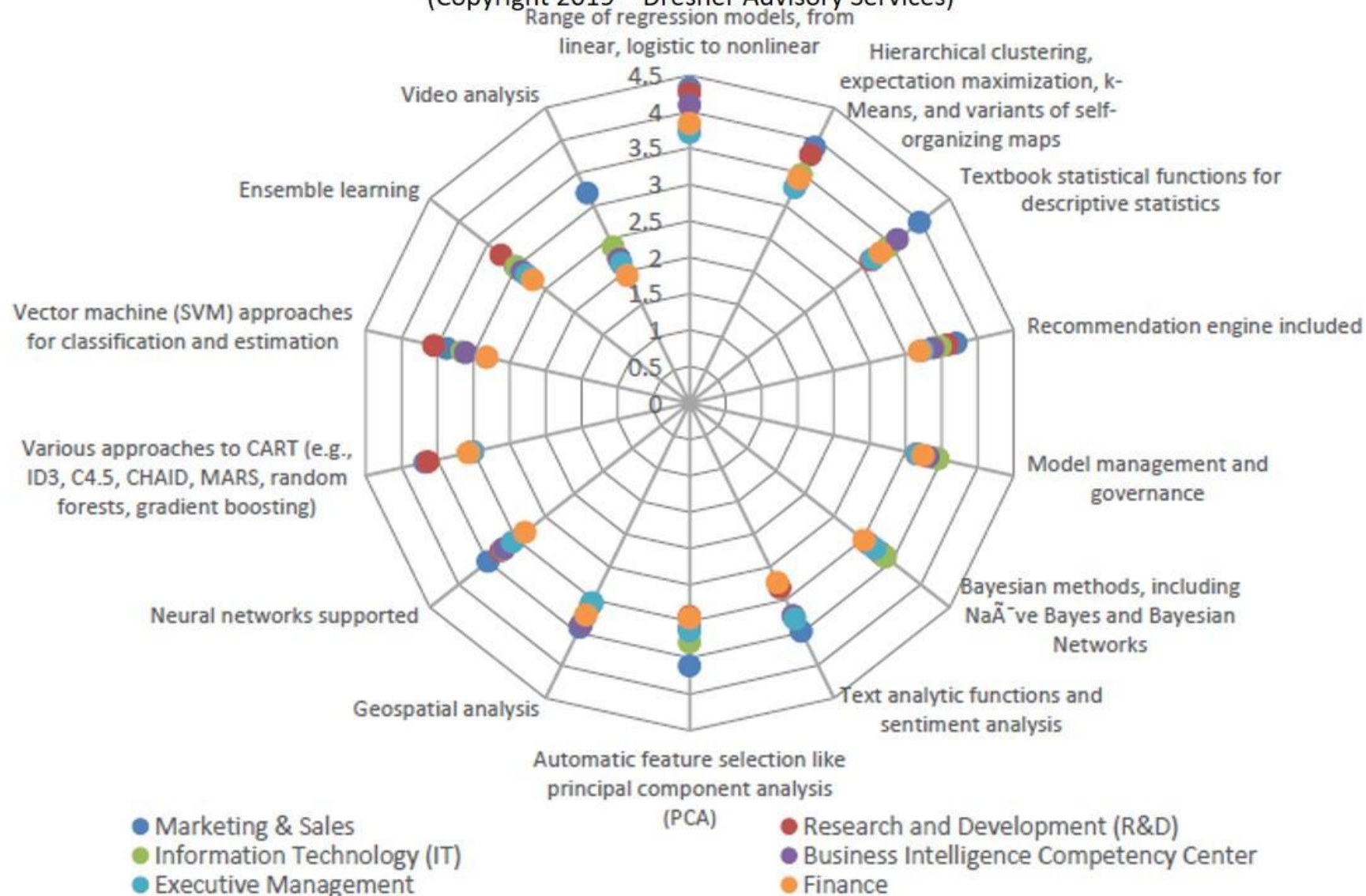
**Louis Columbus**

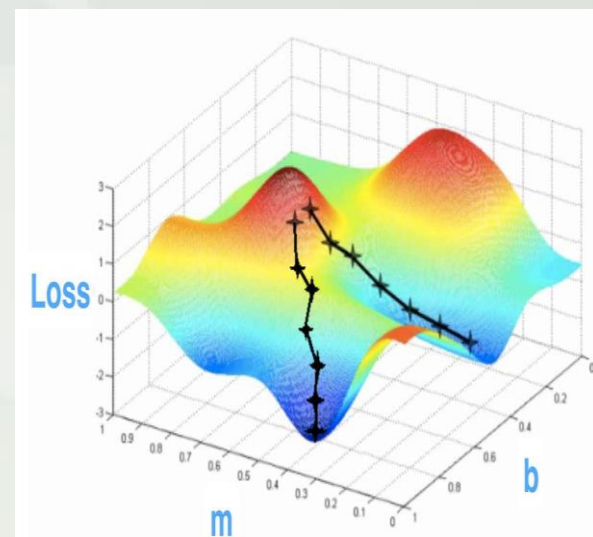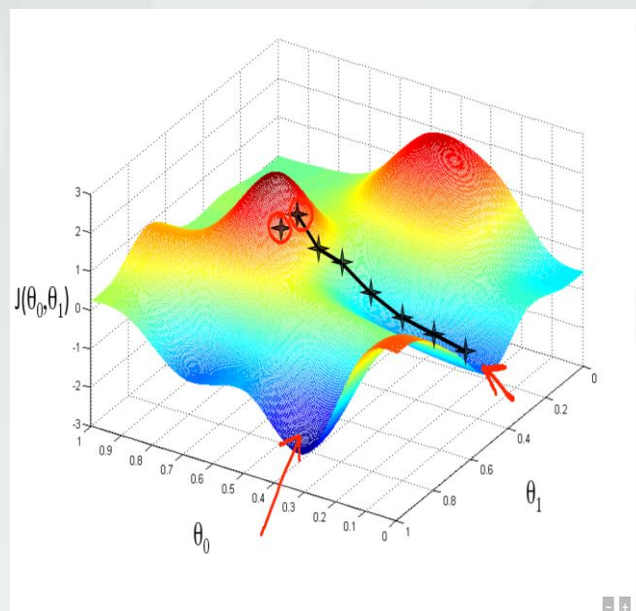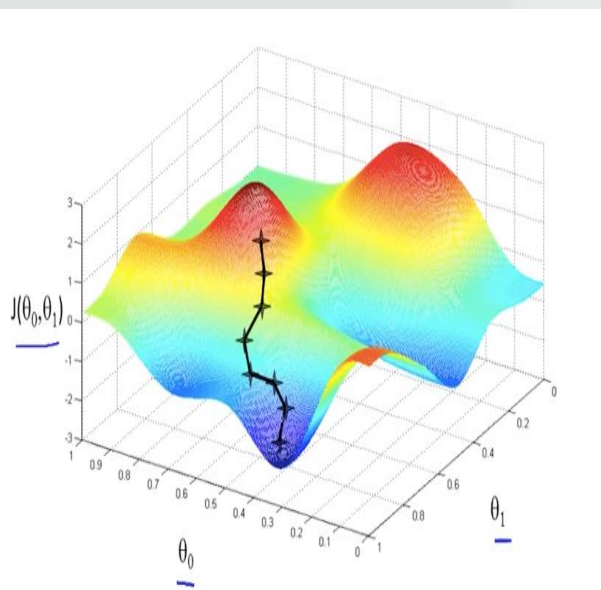# Analytical Features for Data Science and Machine Learning by Function

(Copyright 2019 – Dresner Advisory Services)



- Marketing & Sales
- Information Technology (IT)
- Executive Management
- Research and Development (R&D)
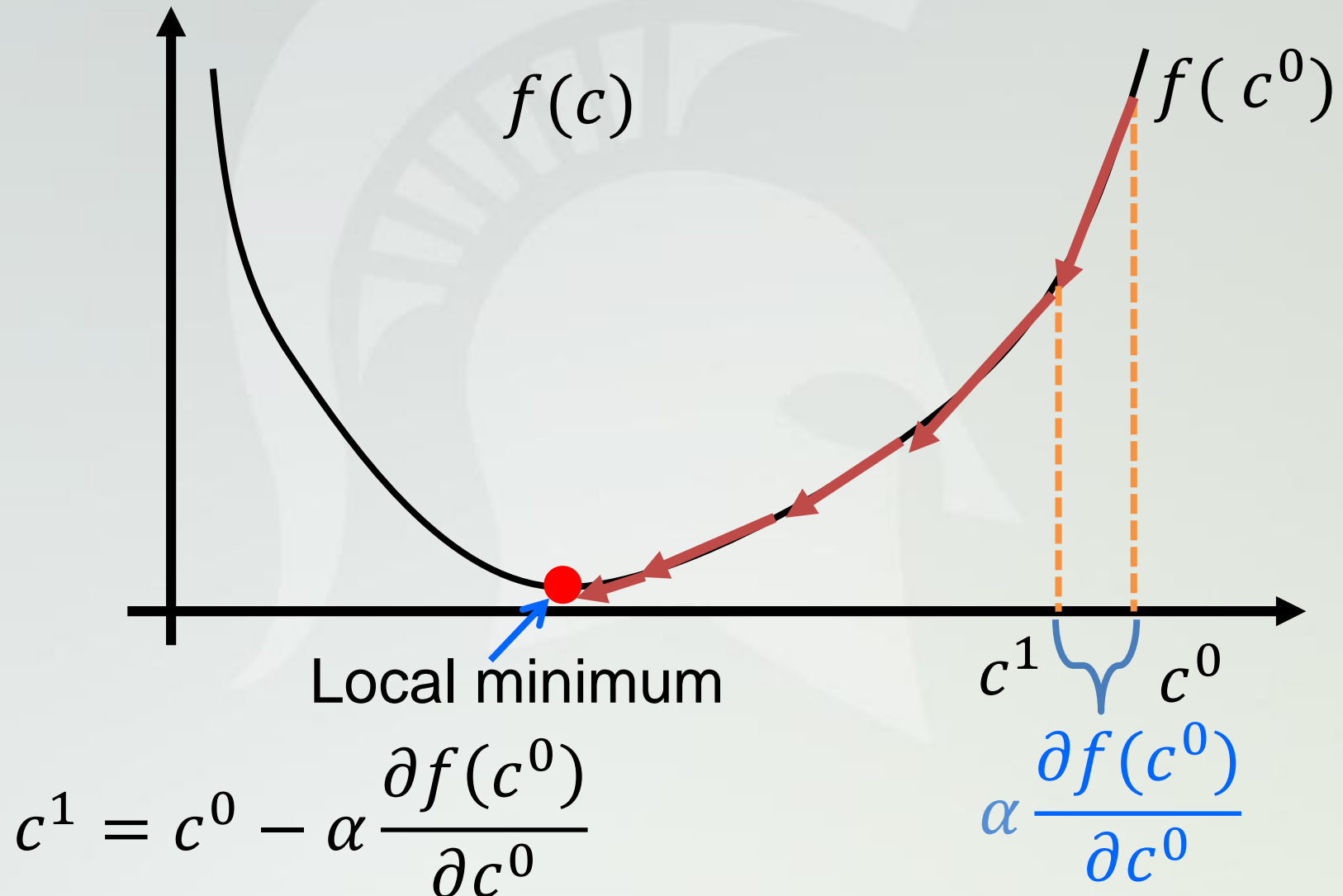- Business Intelligence Competency Center
- Finance

# Introduction

- In general, the loss function has no analytical solutions. We use Gradient Descent (steepest descent or gradient ascent for local maximum).

- Gradient = direction of the steepest ascent

- Find a local minimum of a function

- Often a first-order iterative optimization algorithm

# General Idea

Consider a general function $f(c)$
which can be the loss function of interest



$f(c)$

$f(c^0)$

Local minimum

$c^1$  $c^0$

$\alpha \dfrac{\partial f(c^0)}{\partial c^0}$

$$c^1 = c^0 - \alpha \frac{\partial f(c^0)}{\partial c^0}$$

# **Algorithm**

Find a local minimum of a $C^1$ continuous $f(c)$

- Start with random value $c^0$
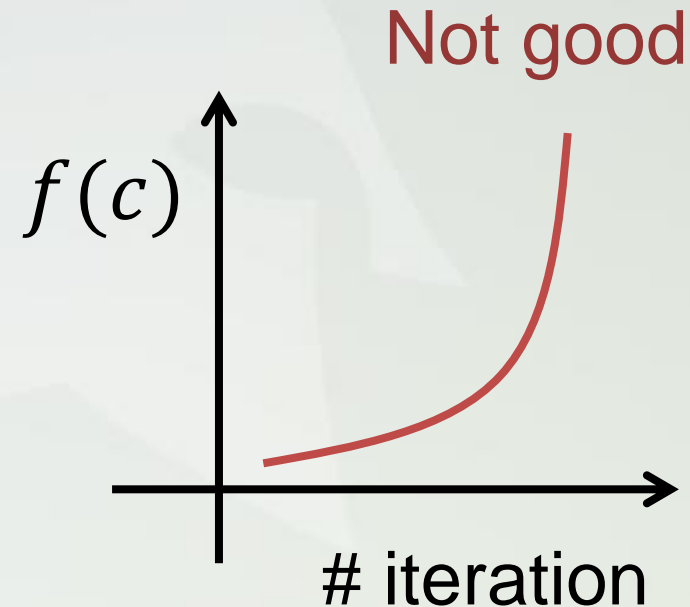
- Update new value:

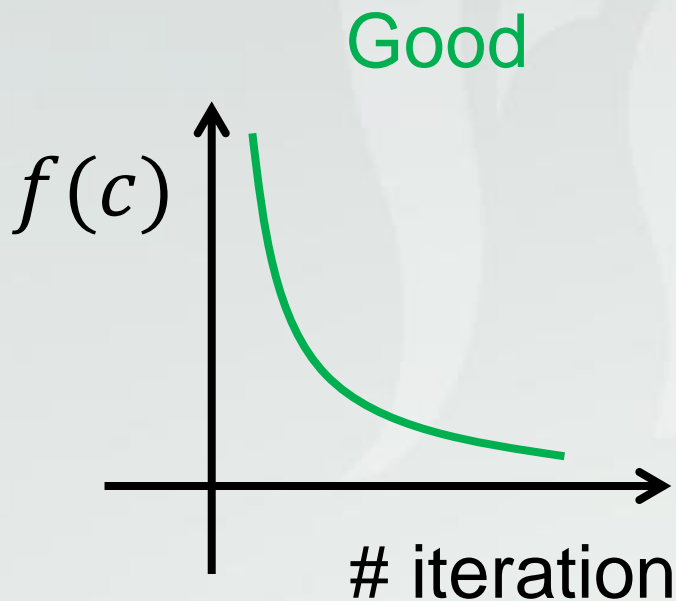$$c^{i+1} = c^i - \alpha \frac{\partial f(c^i)}{\partial c^i}$$

$\alpha$: **learning rate**, very small (like 0.01 or smaller)

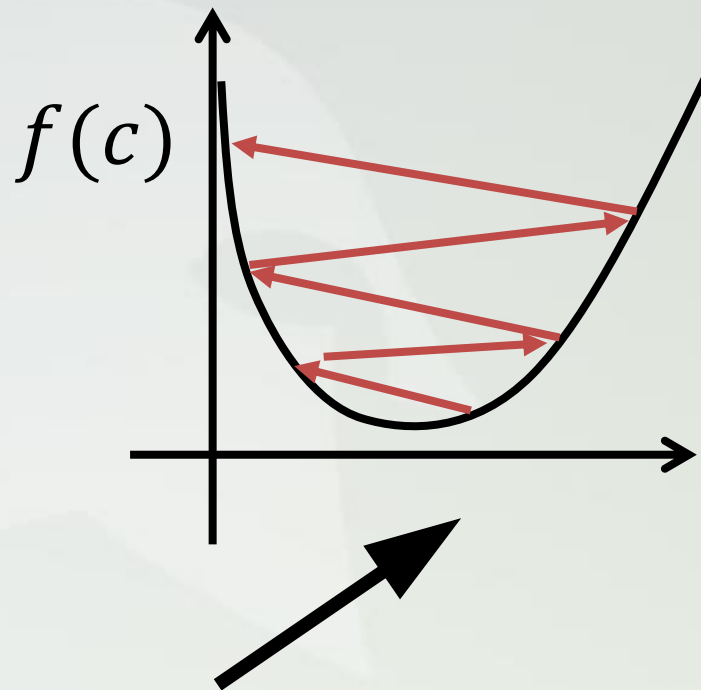- Repeat until $\left\| \frac{\partial f(c^i)}{\partial c^i} \right\| \leq$ tolerance

# Making Sure Gradient Descent Working Correctly

- Function $f(c)$ should decrease after every iteration (monotonically decreases)



Good

$f(c)$

\# iteration

Not good

$f(c)$

\# iteration

# Making Sure Gradient Descent Working Correctly

- Use smaller learning rate $\alpha$

$f(c)$

$f(c)$

Very large learning rate

# Making Sure Gradient Descent Working Correctly

- Feature scaling:

    - Example: assume features for the house price includes number of bedrooms and living area

    - # of bedrooms between 0 and 5

    - But living area between 1 and 5000 $\mathrm{feet}^2$

    - Make all features have the same level of magnitude

# Application for Minimizing Loss Function

- <u>Linear regression</u>: loss function for predictor $p_{\mathbf{c}}(x) = c_0 + c_1 x$ is

$$L(c_0, c_1) = \sum_{i=1}^{M} \left( p\big(x^{(i)}\big) - y^{(i)} \right)^2$$

$$= \sum_{i=1}^{M} \left( c_0 + c_1 x^{(i)} - y^{(i)} \right)^2$$

Use gradient descent to $\min_{c_0, c_1} L(c_0, c_1)$

# Application for Minimizing Loss Function

- Step 1: Assign initial values for $c_0, c_1$:

$$c_0 = 0, c_1 = 1$$

- Step 2: Update the change in values for $c_0, c_1$:

$$c_0 := c_0 - \alpha \frac{\partial}{\partial c_0} L(c_0, c_1)$$

$$:= c_0 - \alpha \sum_{i=1}^{M} 2\left(c_0 + c_1 x^{(i)} - y^{(i)}\right)$$

- **Step 2: (continue)**

$$c_1 := c_1 - \alpha \frac{\partial}{\partial c_1} L(c_0, c_1)$$

$$= c_1 - \alpha \sum_{i=1}^{M} 2x^{(i)}\left(c_0 + c_1 x^{(i)} - y^{(i)}\right)$$

- **Step 3: Repeat Step 2 until it converges**

- <u>Logistic regression</u>: do it similarly

# Other Gradient Descent Methods

- **Stochastic gradient descent (SGD):**
- Herbert Robbins and Sutton Monro (1951)
- Good for large/huge data sets

1) Choose an initial parameter set $c$ and learning rate $\alpha$

2) Randomly shuffle samples in the training set to update $c$

$$c := c - \alpha \frac{\partial}{\partial c} L\big(c, x^{(i)}, y^{(i)}\big), i = 1, 2, .., M$$

(Note: no sum over $i$)

3) Repeat 2) until the convergence is reached.

# Other Gradient Descent Methods

- **SGD with momentum:** accelerate SGD

$$\boldsymbol{v} := \gamma \boldsymbol{v} + \alpha \frac{\partial}{\partial \boldsymbol{c}} L\big(\boldsymbol{c}, \boldsymbol{x}^{(i)}, y^{(i)}\big)$$

$$\boldsymbol{c} := \boldsymbol{c} - \boldsymbol{v}$$

https://distill.pub/2017/momentum/

- Adaptive learning rates are often used.
- If multiple passes are needed, the data can be shuffled for each pass to prevent cycles.

# Other Gradient Descent Methods

- A Method for Stochastic Optimization (Adam) by Kingma & Ba, 2015: An efficiency version of SGD using first and second order momentum, well suited for large data set problems

- Kalman-based Stochastic Gradient Descent: *SIAM Journal on Optimization*. **26** (4): 2620–2648. arXiv:1512.01139

# Other Gradient Descent Methods

## High-order SGD

$$g := \frac{\partial}{\partial \boldsymbol{c}} L\big(\boldsymbol{c}, \boldsymbol{x}^{(i)}, y^{(i)}\big) \qquad \text{(Compute gradient)}$$

$$\boldsymbol{m} := \beta_1 \boldsymbol{m} + (1 - \beta_1)\boldsymbol{g} \qquad \text{(Update 1st order momentum)}$$

$$v := \beta_2 v + (1 - \beta_2)\boldsymbol{g}^2 \qquad \text{(Update 2nd order momentum)}$$

$$\hat{\boldsymbol{m}} := \frac{\boldsymbol{m}}{\beta_1^k} \qquad \text{(Compute corrected-1st order momentum)}$$

$$\hat{v} := \frac{v}{\beta_2^k} \qquad \text{(Compute corrected-2nd order momentum)}$$

$$\boldsymbol{c} := \boldsymbol{c} - \alpha \frac{\hat{\boldsymbol{m}}}{\sqrt{\hat{v}} + \epsilon} \qquad \text{(Update parameters)}$$

# Other Gradient Descent Methods
## Adaptive Gradient Descent

- Barzilai-Bowein method (for $L(\boldsymbol{c})$ convex and $\frac{\partial}{\partial \boldsymbol{c}} L(\boldsymbol{c})$ Lipschitz):

- $\boldsymbol{c}^n = \boldsymbol{c}^{n-1} - \boldsymbol{\alpha}^n \frac{\partial}{\partial \boldsymbol{c}} L(\boldsymbol{c})$

$$\alpha^n = \frac{(\boldsymbol{c}^n - \boldsymbol{c}^{n-1})^T \left[ \frac{\partial}{\partial \boldsymbol{c}} L(\boldsymbol{c}) \Big|_{c=c^n} - \frac{\partial}{\partial \boldsymbol{c}} L(\boldsymbol{c}) \Big|_{c=c^{n-1}} \right]}{\left\| \frac{\partial}{\partial \boldsymbol{c}} L(\boldsymbol{c}) \Big|_{c=c^n} - \frac{\partial}{\partial \boldsymbol{c}} L(\boldsymbol{c}) \Big|_{c=c^{n-1}} \right\|^2}$$

Convex $=>$ the global minimum!

# Discussions

Other potential mathematical approaches:
Explicit Euler, implicit Euler,   Crank–Nicholson, leapfrog, Guass–Legendre Runge–Kutta, Guass–Radau Runge–Kutta, Gauss–Lobatto Runge–Kutta, symplectic Runge–Kutta, Adams–Bashforth, Adams–Moulton, Strong stability preserving, hybrid multistep-multistage methods and adaptive SGD.
(http://users.math.msu.edu/users/wei/paper/p175.pdf)

# Discussions

## Pros and Cons of Gradient Descent

- **Pros**
  - Can be applied for any dimensional space
  - Nonlinear problems
  - Easy to implement

- **Cons:**
  - Local optima problem
  - Slowly to reach the local minimum
  - Cannot be applied for discontinuous functions

# **Discussions**

- **Sample noise (uncertainty in $\{y^{(i)}\}$ )**
- **Parameter linear dependence (in $\{c_i\}$)**
- **Manifold properties:**
  - ➤ **Smoothness -- differentiability**
  - ➤ **Convex/concave**
  - ➤ **Tangent bundle/cotangent bundle**
  - ➤ **Topological structure of the tangent space**
  - ➤ *…*

# Discussions

**Not to be confused with**

- **Method of steepest descent (for integrals)**
- **Conjugated gradient method**