

## **Deep Learning for Modeling: Concepts, Tools, and Techniques**

### **Week 6: Momentum-based Gradient Descent and Learning Rates**

Li Shuo-Hui

## Gradient-based optimization: revisit

$$\min \mathcal{L}(\mathbf{x})$$

**Update:**  $\mathbf{x}_{k+1} = \mathbf{x}_k + \eta_k \mathbf{p}_k$

$\eta_k$   $\rightarrow$  *Learning rate/step size*

$\mathbf{p}_k = -\mathbf{B}_k^{-1} \frac{\partial \mathcal{L}}{\partial \mathbf{x}_k}$

$\mathbf{B}_k^{-1}$   $\rightarrow$  *Local geometry*

$\mathbf{p}_k$   $\rightarrow$  *Optim direction*

## Gradient descent: revisit

### Gradient descent

*Repeat*

$$\boldsymbol{\theta}' = \boldsymbol{\theta} - \eta \frac{\partial \mathcal{L}_{\boldsymbol{\theta}}(\mathbf{x})}{\partial \boldsymbol{\theta}}$$

*Till  $\mathcal{L}_{\boldsymbol{\theta}}(\mathbf{x})$  is small enough*

$$\eta_k = \eta$$

$$\mathbf{p}_k = -\frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}}$$

$$\mathbf{B}_k^{-1} = \mathbf{I}$$

## Newton's method: revisit

## Newton's method

Solve  $\operatorname{argmin}_{\mathbf{x}} f(\mathbf{x})$

$$\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = 0 \longrightarrow \frac{\partial f}{\partial \mathbf{x}} + \delta \mathbf{x} \frac{\partial^2 f}{\partial \mathbf{x}^2} = 0$$

Repeat  $\mathbf{x}' = \mathbf{x} - \delta \mathbf{x}$  till converge

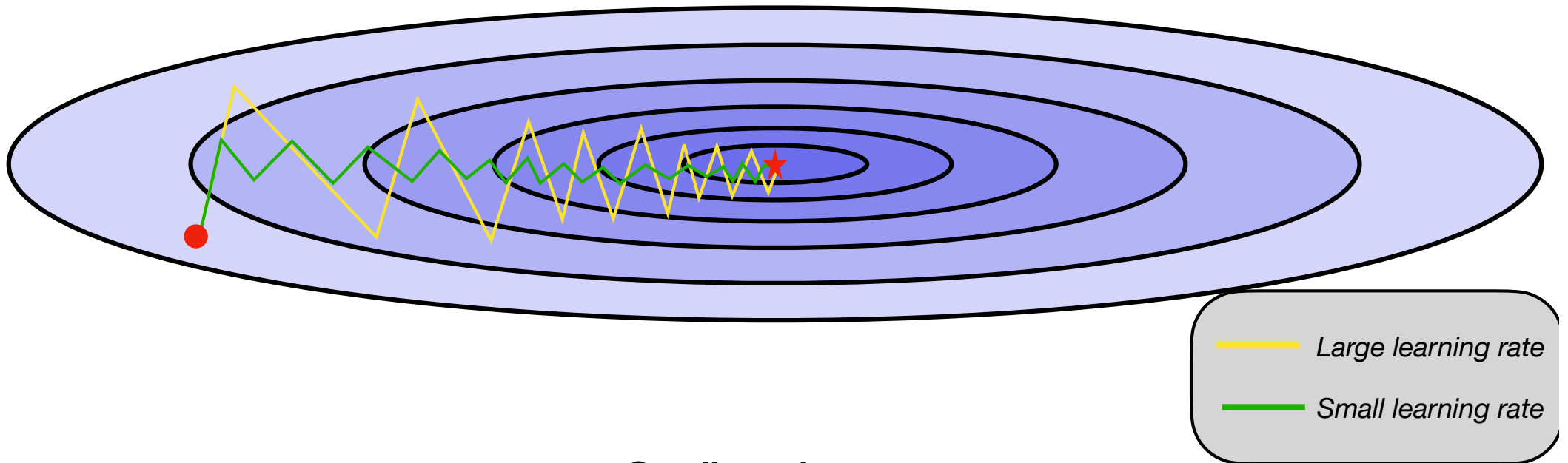
$$\eta_k = 1$$

$$\mathbf{p}_k = -\left(\frac{\partial^2 f}{\partial \mathbf{x}^2}\right)^{-1} \frac{\partial f}{\partial \mathbf{x}}$$

$$\mathbf{B}_k^{-1} = \left(\frac{\partial^2 f}{\partial \mathbf{x}^2}\right)^{-1}$$

## Gradient-based optimization: problem one

### The local geometry

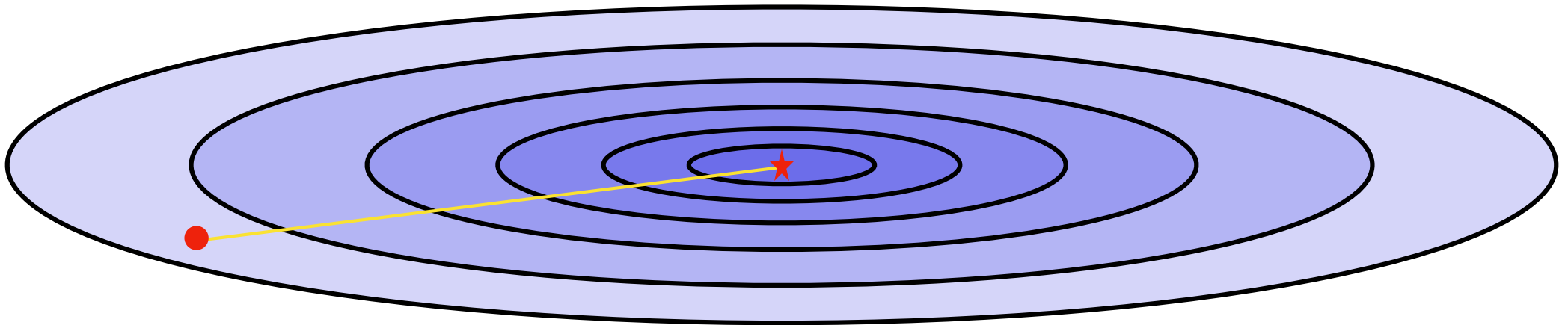


Gradient descent

$$\mathbf{B}_k^{-1} = \mathbf{I}$$

## Gradient-based optimization: problem one

### The local geometry

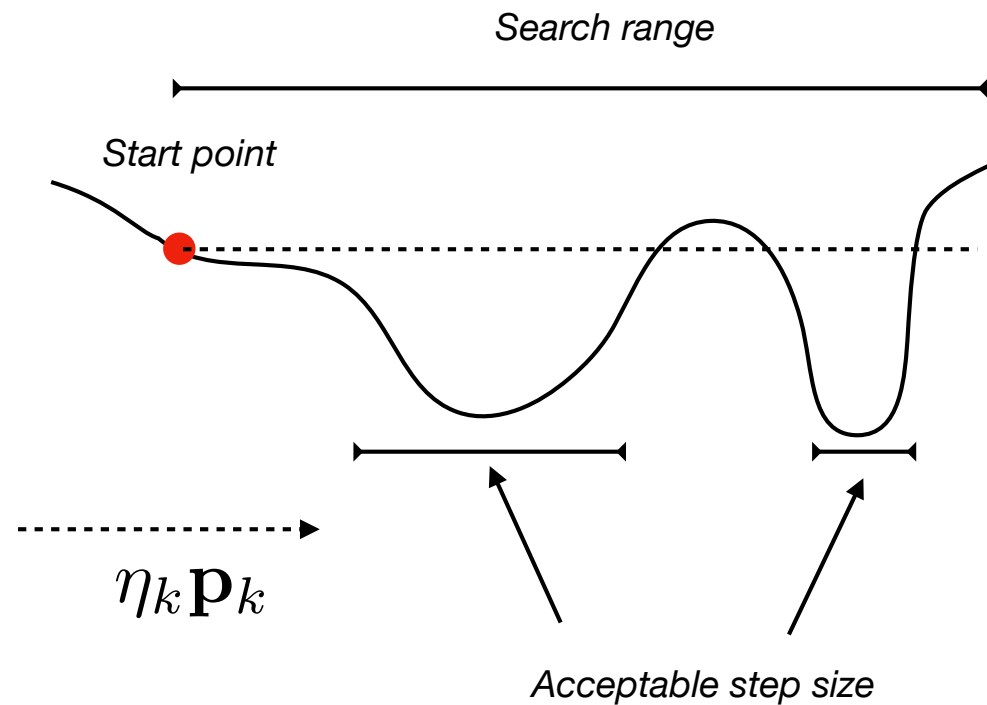


### Newton's method

$$\mathbf{B}_k^{-1} = \left( \frac{\partial^2 f}{\partial \mathbf{x}^2} \right)^{-1} \quad \text{Balance update across directions!}$$

## Gradient-based optimization: problem two

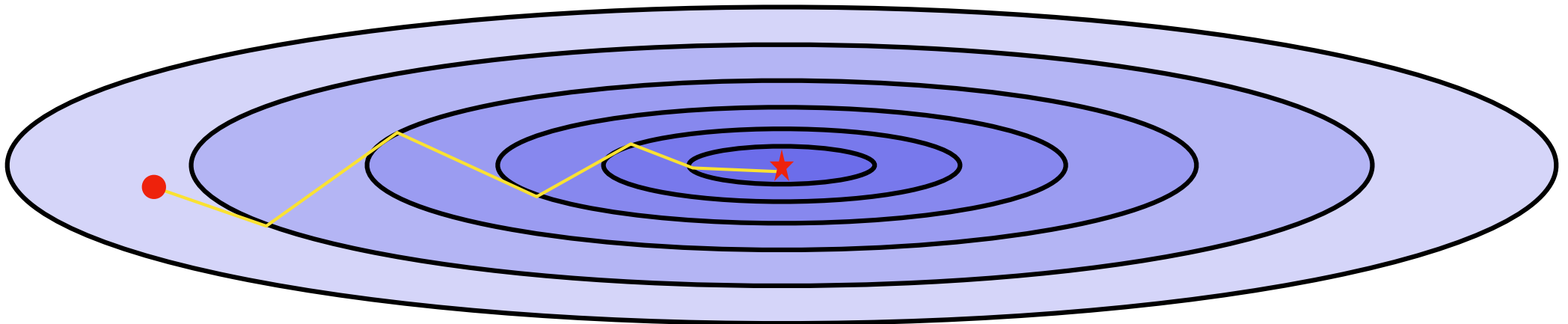
### The step size selection



### The line search algorithms

## Gradient-based optimization: problem one

### The step size selection









**Steepest Descent:**  
gradient descent with exact line search

$$\mathbf{B}_k^{-1} = \mathbf{I} \quad \text{“Zig-zag track” tangents to level set}$$



## Gradient-based optimization

### Comparison

	Line search step size	Local geometry acknowledge
Newton /w line search		
Steepest descent		
SGD		

*These can be fixed with momentum and changing learning rate*

## Momentum-based gradient descent

### Simple momentum

$$\mathbf{p}_k = -\rho \mathbf{p}_{k-1} - \frac{\partial f}{\partial \mathbf{x}_k}$$

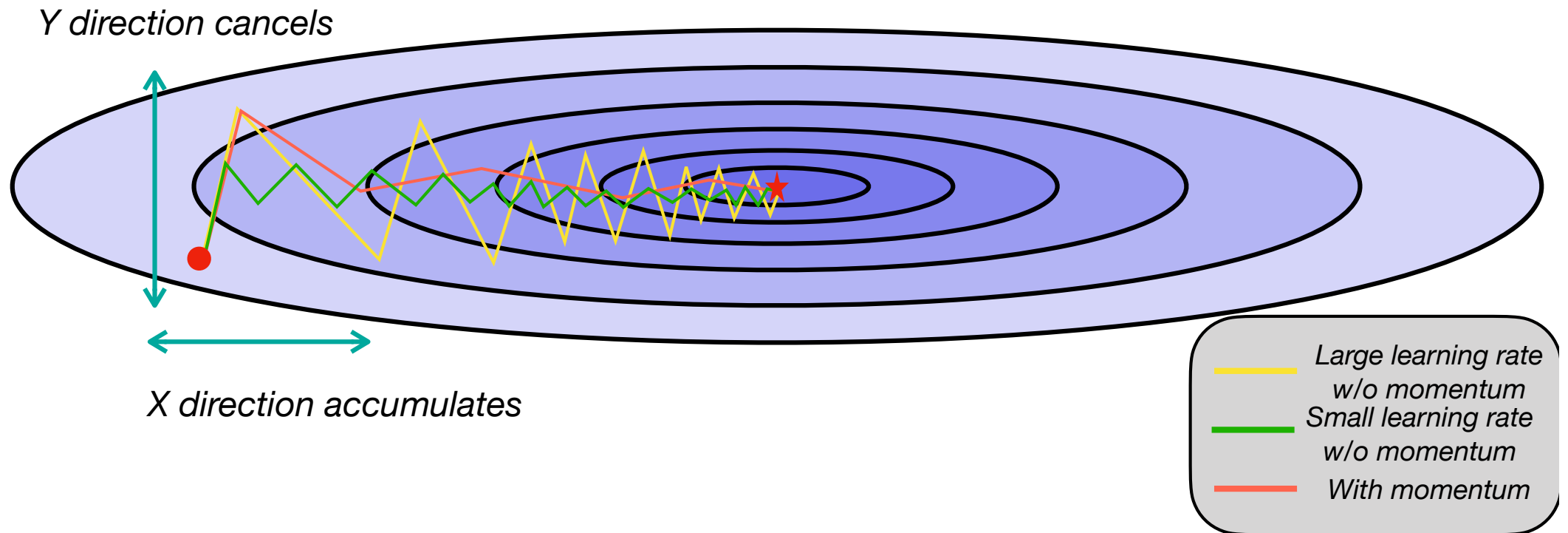
$$\mathbf{x}_{k+1} = \mathbf{x}_k + \eta \mathbf{p}_k$$

*Momentum term*

*Momentum decay factor  
Or “friction”*

## Momentum-based gradient descent

$$\mathbf{p}_k = -\rho \mathbf{p}_{k-1} - \frac{\partial f}{\partial \mathbf{x}_k}$$



## Changing learning rates

### Decaying learning rate

*Exponential*

$$\eta_k = \eta_0 \cdot \gamma^{\max(0, \lceil \frac{k-k_0}{s} \rceil)}$$

*1/t scheme*

$$\eta_k = \frac{\eta_0}{1 + \max(0, \lceil \frac{k-k_0}{s} \rceil)}$$

⋮

**Decaying learning rate to adapt finer optimization at the later stage**

## Changing learning rates

### Adaptive learning rate

*Adagrad*

$$\eta_{k,i} = \frac{\eta_0}{\sqrt{\sum_j^k p_{j,i}^2}}$$

$$\begin{pmatrix} x_{k+1,1} \\ x_{k+1,2} \\ \vdots \\ x_{k+1,n} \end{pmatrix} = \begin{pmatrix} x_{k,1} \\ x_{k,2} \\ \vdots \\ x_{k,n} \end{pmatrix} + \begin{pmatrix} \frac{\eta_0}{\sqrt{\sum_j^k p_{j,1}^2}} \\ \frac{\eta_0}{\sqrt{\sum_j^k p_{j,2}^2}} \\ \vdots \\ \frac{\eta_0}{\sqrt{\sum_j^k p_{j,n}^2}} \end{pmatrix} \odot \begin{pmatrix} p_{k,1} \\ p_{k,2} \\ \vdots \\ p_{k,n} \end{pmatrix}$$

**Adaptively tune the learning rate per parameter using historical updates**

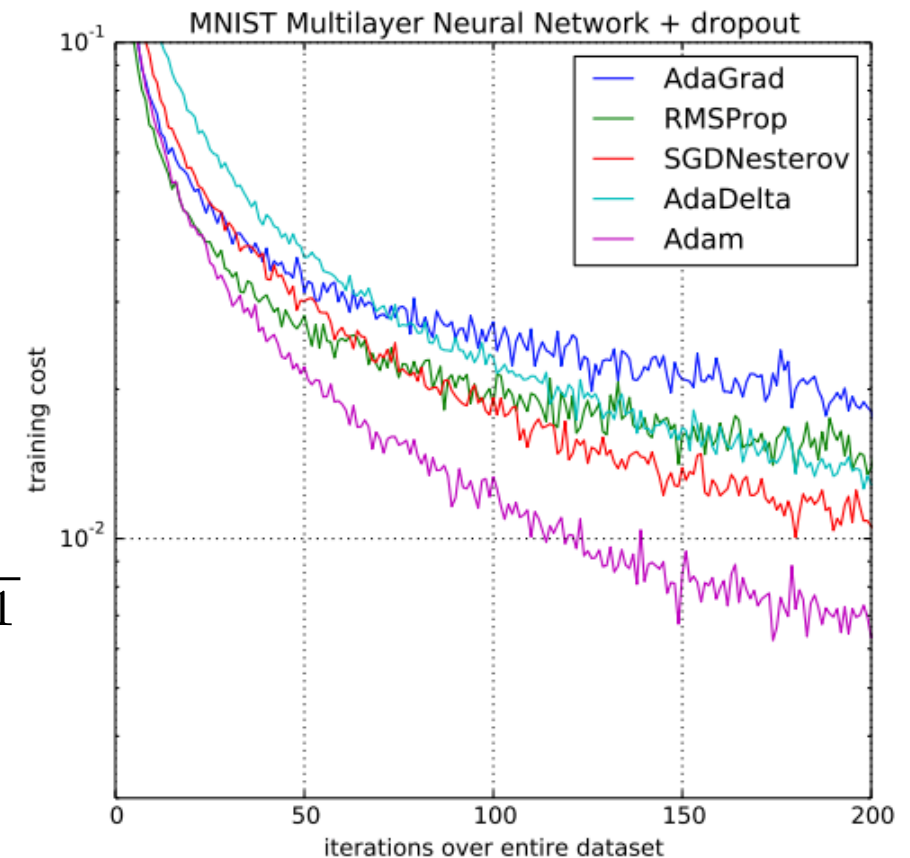
## Adam optimizer

$$m_{k+1} = \beta_1 m_k + (1 - \beta_1) \frac{\partial \mathcal{L}}{\partial \theta_k}$$

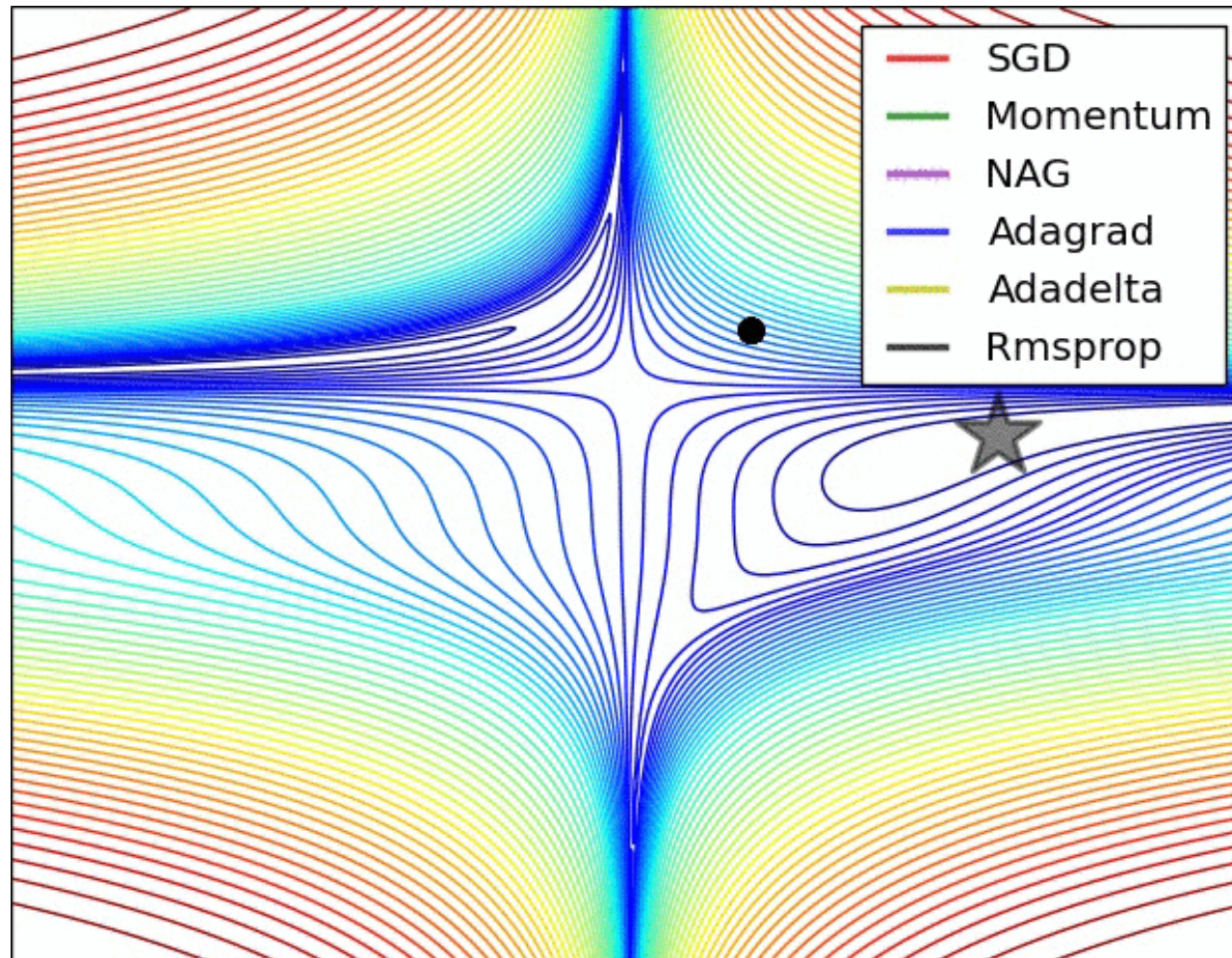
$$v_{k+1} = \beta_2 v_k + (1 - \beta_2) \frac{\partial \mathcal{L}}{\partial \theta_k}^2$$

$$\hat{m}_{k+1} = \frac{m_{k+1}}{1 - \beta_1^{k+1}} \quad \hat{v}_{k+1} = \frac{v_{k+1}}{1 - \beta_2^{k+1}}$$

$$\theta_{k+1} = \theta_k - \eta \frac{\hat{m}_{k+1}}{\sqrt{\hat{v}_{k+1}}}$$



## Some intuitions



# **Deep Learning for Modeling: Concepts, Tools, and Techniques**

## Week 6 Tutorial

Li Shuo-Hui