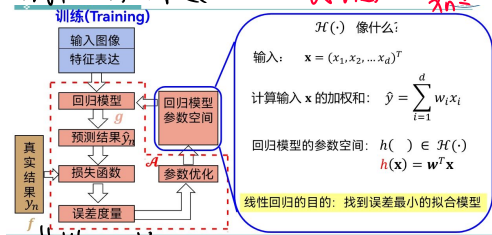
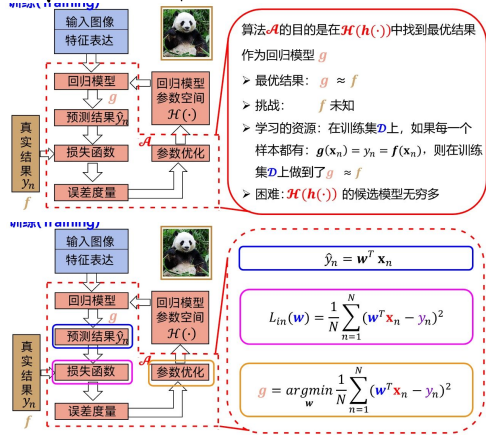


### 3.1 线性回归问题



### 3.2 线性回归算法



### 有矩阵/向量开式 $L_{in}(\mathbf{w})$

$L_{in}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n - y_n)^2 = \frac{1}{N} \left\| \begin{bmatrix} \mathbf{x}_1^T \mathbf{w} - y_1 \\ \vdots \\ \mathbf{x}_N^T \mathbf{w} - y_N \end{bmatrix} \right\|^2 = \frac{1}{N} \left\| \begin{bmatrix} -\mathbf{x}_1^T & -1 \\ \vdots & \vdots \\ -\mathbf{x}_N^T & -1 \end{bmatrix} \mathbf{w} - \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} \right\|^2$

$\mathbf{X}: N \times (d+1)$   
 $\mathbf{w}: (d+1) \times 1$   
 $\mathbf{Y}: N \times 1$

$L_{in}(\mathbf{w}) = \frac{1}{N} \|\mathbf{X}\mathbf{w} - \mathbf{Y}\|^2$

求最优解  $\min_{\mathbf{w}} L_{in}(\mathbf{w}) = \frac{1}{N} \|\mathbf{X}\mathbf{w} - \mathbf{Y}\|^2 \Rightarrow \nabla L_{in}(\mathbf{w}) = 0$  有最佳解

连续、可微、凸函数

$\nabla L_{in}(\mathbf{w}) = \begin{bmatrix} \frac{\partial L_{in}(\mathbf{w})}{\partial w_0} \\ \frac{\partial L_{in}(\mathbf{w})}{\partial \mathbf{w}} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

$\nabla L_{in}(\mathbf{w})$  求解:

$L_{in}(\mathbf{w}) = \frac{1}{N} \|\mathbf{X}\mathbf{w} - \mathbf{Y}\|^2 = \frac{1}{N} (\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{Y} + \mathbf{Y}^T \mathbf{Y}) = \frac{1}{N} (\mathbf{w}^T \mathbf{A} \mathbf{w} - 2\mathbf{w}^T \mathbf{b} + c)$

当  $\mathbf{w}$  是单变量时:  $L_{in}(w) = \frac{1}{N} (aw^2 - 2bw + c)$   
 $\nabla L_{in}(w) = \frac{1}{N} (2aw - 2b)$

当  $\mathbf{w}$  是向量时:  $L_{in}(\mathbf{w}) = \frac{1}{N} (\mathbf{w}^T \mathbf{A} \mathbf{w} - 2\mathbf{w}^T \mathbf{b} + c)$   
 $\nabla L_{in}(\mathbf{w}) = \frac{1}{N} (2\mathbf{A}\mathbf{w} - 2\mathbf{b})$

$\nabla L_{in}(\mathbf{w}) = \frac{2}{N} (\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{Y})$

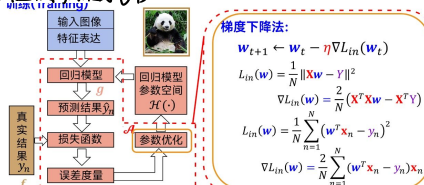
### 3.3 GD

#### 感知器算法的GD

- 对样本的特征向量  $\mathbf{x}$  和权重向量  $\mathbf{w}$  推广化
- 初始化权重向量  $\mathbf{w}_0$  (例如:  $\mathbf{w}_0 = 0$ )
- for  $t = 0, 1, 2, \dots$  ( $t$  代表迭代次数)
  - 对某些样本  $n$ , 通过下式对权重向量  $\mathbf{w}_t$  进行更新:  
 $\mathbf{w}_{t+1} = \mathbf{w}_t + \frac{1}{\eta} \cdot (\text{sign}(\mathbf{w}_t^T \mathbf{x}_{n(t)}) - y_n) \mathbf{x}_{n(t)}$
- 直到满足停止条件, 此时的  $\mathbf{w}_{t+1}$  作为学到的  $\mathbf{g}$

算法可理解成通过选择  $(\eta, \mathbf{v})$ , 以及确定“停止条件”的找到最佳解的迭代优化过程

#### 线性回归的GD



- 初始化权重向量  $\mathbf{w}_0$
- for  $t = 0, 1, 2, \dots$  ( $t$  代表迭代次数)
  - 计算梯度:  $\nabla L_{in}(\mathbf{w}) = \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n - y_n) \mathbf{x}_n$
  - 对权重向量  $\mathbf{w}_t$  进行更新:  $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta \nabla L_{in}(\mathbf{w}_t)$
- 直到  $\nabla L_{in}(\mathbf{w}) = 0$ , 或者迭代足够多次
- 返回最终的  $\mathbf{w}_{t+1}$  作为学到的  $\mathbf{g}$

#### 据梯度调整学习率引出改良优化器

① Adagrad

$\sigma_{i,t+1} \leftarrow \sigma_{i,t} + \frac{\eta}{\sigma_{i,t}} \frac{\partial L_{in}}{\partial w_{i,t}}$

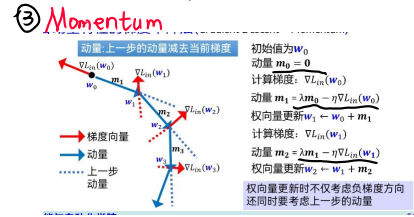
② RMS prop

$\sigma_{i,t+1} \leftarrow \sigma_{i,t} + \frac{\eta}{\sigma_{i,t}} \frac{\partial L_{in}}{\partial w_{i,t}}$

$\sigma_{i,t} = \sqrt{\frac{1}{t+1} \sum_{\tau=0}^t \left( \frac{\partial L_{in}}{\partial w_{i,\tau}} \right)^2}$

通过  $\alpha$  的取值, 使得当前的梯度影响更大, 而以往的梯度影响较小

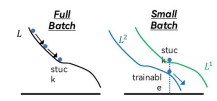
#### 逃离鞍点/局部极值引出改良优化器



#### ④ 结合两者 Adam

**Algorithm 1: Adam**, our proposed algorithm for stochastic optimization. See section 2 for details, and for a slightly more efficient (but less clear) order of computation.  $g_t^i$  indicates the elementwise square  $g_t \odot g_t$ . Good default settings for the tested machine learning problems are  $\alpha = 0.001$ ,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$  and  $\epsilon = 10^{-8}$ . All operations on vectors are element-wise. With  $\beta_1^t$  and  $\beta_2^t$  we denote  $\beta_1$  and  $\beta_2$  to the power  $t$ .

**Require:**  $\alpha$ : Stepsize  
**Require:**  $\beta_1, \beta_2 \in [0, 1]$ : Exponential decay rates for the moment estimates  
**Require:**  $f(\theta)$ : Stochastic objective function with parameters  $\theta$   
**Require:**  $\theta_0$ : Initial parameter vector  
 $\mathbf{m}_0 \leftarrow 0$  (Initialize 1<sup>st</sup> moment vector)  $\rightarrow$  for momentum  
 $\mathbf{v}_0 \leftarrow 0$  (Initialize 2<sup>nd</sup> moment vector)  $\rightarrow$  for RMSprop  
 $t \leftarrow 0$  (Initialize timestep)  
**while**  $\theta_t$  not converged **do**  
   $t \leftarrow t + 1$   
   $g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$  (Get gradients w.r.t. stochastic objective at timestep  $t$ )  
   $\mathbf{m}_t \leftarrow \beta_1 \cdot \mathbf{m}_{t-1} + (1 - \beta_1) \cdot g_t$  (Update biased first moment estimate)  
   $\mathbf{v}_t \leftarrow \beta_2 \cdot \mathbf{v}_{t-1} + (1 - \beta_2) \cdot g_t^2$  (Update biased second raw moment estimate)  
   $\hat{\mathbf{m}}_t \leftarrow \mathbf{m}_t / (1 - \beta_1^t)$  (Compute bias-corrected first moment estimate)  
   $\hat{\mathbf{v}}_t \leftarrow \mathbf{v}_t / (1 - \beta_2^t)$  (Compute bias-corrected second raw moment estimate)  
   $\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \hat{\mathbf{m}}_t / (\sqrt{\hat{\mathbf{v}}_t} + \epsilon)$  (Update parameters)  
**end while**  
**return**  $\theta_t$  (Resulting parameters)



#### 批量用于梯度下降

批量  $\uparrow \Rightarrow$  Batch 速度  $\downarrow$  epoch 速度  $\uparrow$  性能  $\downarrow$

- 初始化权重向量  $\mathbf{w}_0$
  - for  $t = 0, 1, 2, \dots$  ( $t$  代表迭代次数)
    - 计算梯度:  $\nabla L_{in}(\mathbf{w}) = \frac{2}{N} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}_n - y_n) \mathbf{x}_n$
    - 对权重向量  $\mathbf{w}_t$  进行更新:  $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta \nabla L_{in}(\mathbf{w}_t)$
  - 直到  $\nabla L_{in}(\mathbf{w}) = 0$ , 或者迭代足够多次
  - 返回最终的  $\mathbf{w}_{t+1}$  作为学到的  $\mathbf{g}$
- 随机梯度下降法 (Stochastic Gradient Descent) (SGD)

	批量小	批量大
一次更新需要的速度 (无并行处理)	Faster	Slower
一次更新需要的速度 (有并行处理)	Same	Same (not too large)
一个 epoch 花费的时间	Slower	Faster
梯度的特点	Noisy	Stable
优化性能	Better	Worse
泛化性能	Better	Worse