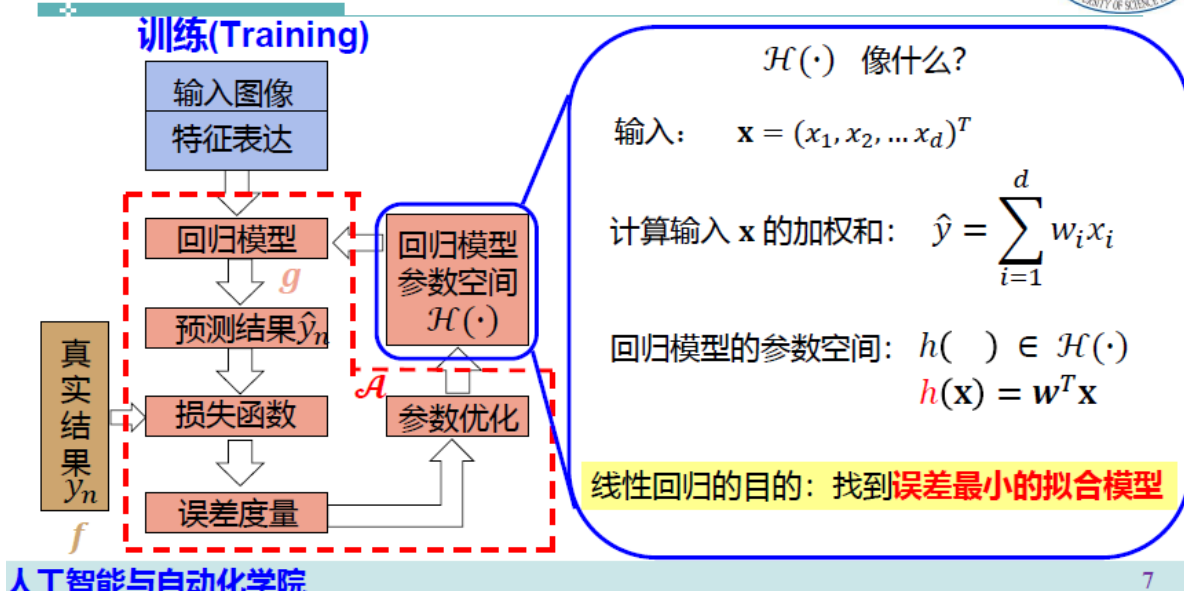# 模式识别 U3 线性回归 Linear Regression

## 课程内容

- 3.1 线性回归问题
- 3.2 线性回归算法
- 3.3 梯度下降法

## 线性回归问题

机器学习的过程其实是一个找**最拟合函数**的过程，通过不断的训练，我们最终得到一个函数映射，给定函数（网络）一个输入，函数（网络）会给出相应的输出

**若输出的是一个数值（scatter），我们就将这类机器学习问题称为回归Regression**

**训练(Training)**

输入图像 特征表达

回归模型 $g$

预测结果 $\hat{y}_n$

真实结果 $y_n$ $f$

损失函数

误差度量

回归模型 参数空间 $\mathcal{H}(\cdot)$

参数优化

$\mathcal{A}$

$\mathcal{H}(\cdot)$ 像什么?

输入: $\mathbf{x} = (x_1, x_2, \ldots x_d)^T$

计算输入 $\mathbf{x}$ 的加权和: $\hat{y} = \sum_{i=1}^{d} w_i x_i$

回归模型的参数空间: $h(\ ) \in \mathcal{H}(\cdot)$
$$h(\mathbf{x}) = \boldsymbol{w}^T \mathbf{x}$$

线性回归的目的: 找到**误差最小的拟合模型**

# 线性回归算法

## 模型构建

**训练(Training)**

输入图像 特征表达

回归模型 $g$

预测结果 $\hat{y}_n$

真实结果 $y_n$ $f$

损失函数

误差度量

回归模型 参数空间 $\mathcal{H}(\cdot)$

参数优化

$\mathcal{A}$

算法 $\mathcal{A}$ 的目的是在 $\mathcal{H}(h(\cdot))$ 中找到最优结果作为回归模型 $g$

➢ 最优结果: $g \approx f$

➢ 挑战: $f$ 未知

➢ 学习的资源: 在训练集 $\mathcal{D}$ 上,如果每一个样本都有: $g(\mathbf{x}_n) = y_n = f(\mathbf{x}_n)$,则在训练集 $\mathcal{D}$ 上做到了 $g \approx f$

➢ 困难: $\mathcal{H}(h(\cdot))$ 的候选模型无穷多

为达到回归目的,我们度量模型输出结果时不再仅仅关注输出的符号

而是关注模型输出的数值

**由此，我们择取 "平方误差函数"作为我们的损失函数，来度量我们的模型学习效果**

$$\mathcal{L} = (\hat{y_n} - y_n)^2$$

$$\mathcal{L}_{in} = \frac{1}{N} \sum_{n=1}^{N} (\hat{y_n} - y_n)^2$$

$\mathcal{L}_{in}$计算训练样本集所有样本产生的平均损失

$$\mathcal{L}_{in}(h) = \frac{1}{N} \sum_{n=1}^{N} (h(\mathbf{x}_n) - y_n)^2$$

$h(\mathbf{x}_n)$是回归模型的结果

**在线性回归模型中，** $\hat{y_n} = h(\mathbf{x}_n) = \mathbf{w}^T \mathbf{x}_n$

$$\mathcal{L}_{in}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^{N} (\mathbf{w}^T \mathbf{x}_n - y_n)^2$$

则

$$g = \arg \min_{\mathbf{w}} \frac{1}{N} \sum_{n=1}^{N} (\mathbf{w}^T \mathbf{x}_n - y_n)^2$$

## 向量/矩阵形式

$$L_{in}(\boldsymbol{w}) = \frac{1}{N}\sum_{n=1}^{N}(\boldsymbol{w}^T\mathbf{x}_n - y_n)^2 = \frac{1}{N}\sum_{n=1}^{N}(\mathbf{x}_n^T\boldsymbol{w} - y_n)^2$$

$$= \frac{1}{N}\left\|\begin{bmatrix}\mathbf{x}_1^T\boldsymbol{w} - y_1\\\mathbf{x}_2^T\boldsymbol{w} - y_2\\\cdots\\\mathbf{x}_N^T\boldsymbol{w} - y_N\end{bmatrix}\right\|^2 = \frac{1}{N}\left\|\begin{bmatrix}--\mathbf{x}_1^T--\\--\mathbf{x}_2^T--\\\cdots\\--\mathbf{x}_N^T--\end{bmatrix}\boldsymbol{w} - \begin{bmatrix}y_1\\y_2\\\cdots\\y_N\end{bmatrix}\right\|^2$$

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

$\boldsymbol{w}$: $(d+1) \times 1$

$Y$: $N \times 1$

$$= \frac{1}{N}\|\mathbf{X}\boldsymbol{w} - Y\|^2$$

**求最佳解：** $\min\limits_{\boldsymbol{w}} L_{in}(\boldsymbol{w}) = \frac{1}{N}\|\mathbf{X}\boldsymbol{w} - Y\|^2$

$L_{in}$ 曲线具有连续、可微、凸函数的特点



$$\nabla L_{in}(\boldsymbol{w}) = \begin{bmatrix}\dfrac{\partial L_{in}(\boldsymbol{w})}{\partial w_0}\\\dfrac{\partial L_{in}(\boldsymbol{w})}{\partial w_1}\\\vdots\\\dfrac{\partial L_{in}(\boldsymbol{w})}{\partial w_d}\end{bmatrix} = \begin{bmatrix}0\\0\\\vdots\\0\end{bmatrix}$$

$\nabla L_{in}(\boldsymbol{w}) = \mathbf{0}$时求得最佳解$\boldsymbol{w}^*$

# 广义逆解法——线性回归的解析解

$\nabla\mathcal{L}_{in}(\mathbf{w}) = 0$时，求得最佳解$\mathbf{w}^*$

$\nabla L_{in}(\boldsymbol{w})$**求解：**

$$L_{in}(\boldsymbol{w}) = \frac{1}{N}\|\mathbf{X}\boldsymbol{w} - Y\|^2 = \frac{1}{N}(\boldsymbol{w}^T\mathbf{X}^T\mathbf{X}\boldsymbol{w} - 2\boldsymbol{w}^T\mathbf{X}^TY + Y^TY) = \frac{1}{N}(\boldsymbol{w}^TA\boldsymbol{w} - 2\boldsymbol{w}^Tb + c)$$

当$w$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$\boldsymbol{w}$是向量时

$$L_{in}(\boldsymbol{w}) = \frac{1}{N}(\boldsymbol{w}^TA\boldsymbol{w} - 2\boldsymbol{w}^Tb + c)$$

$$\nabla L_{in}(\boldsymbol{w}) = \frac{1}{N}(2A\boldsymbol{w} - 2b)$$

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

$$\nabla L_{in}(\boldsymbol{w}) = \frac{2}{N}(\mathbf{X}^T\mathbf{X}\boldsymbol{w} - \mathbf{X}^T Y) = \mathbf{0}$$

$$\boldsymbol{w}^* = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T Y$$

$$g = \boldsymbol{w}^* = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T Y = \mathbf{X}^\dagger Y$$

$$\mathbf{X}^\dagger = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T, \text{ 称为广义逆}$$

**线性回归算法**

- 对训练样本集 $\mathcal{D}$ 构造输入特征向量矩阵 $\mathbf{X}$ 和输出向量 $Y$

$$\mathbf{X} = \begin{bmatrix} --\mathbf{x}_1^T-- \\ --\mathbf{x}_2^T-- \\ \vdots \\ --\mathbf{x}_N^T-- \end{bmatrix}, \qquad Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

- 计算广义逆: $\mathbf{X}^\dagger = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$

- 计算回归值: $g = \boldsymbol{w}^* = \mathbf{X}^\dagger Y$

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

$Y$: $N \times 1$

$\mathbf{X}^\dagger$: $(d+1) \times N$

$\boldsymbol{w}$: $(d+1) \times 1$

# 梯度下降法 Gradient Descent

## 各种 GD算法



训练(Training)

输入图像 特征表达

回归模型 $g$

预测结果 $\hat{y}_n$

损失函数

误差度量

真实结果 $y_n$  $f$

回归模型参数空间 $\mathcal{H}(\cdot)$

$\mathcal{A}$ 参数优化

$$\nabla L_{in}(\boldsymbol{w}) = \frac{2}{N}(\mathbf{X}^T\mathbf{X}\boldsymbol{w} - \mathbf{X}^T Y) = \mathbf{0}$$

$$\boldsymbol{w}^* = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T Y$$

$$g = \boldsymbol{w}^* = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T Y = \mathbf{X}^\dagger Y$$

$$\mathbf{X}^\dagger = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T, \text{ 称为广义逆}$$

当 $N = 10000$, $d = 9999$, $\mathbf{X}$ 是一个巨大的矩阵, $N$, $d$ 可以更大, 如何计算 $\mathbf{X}^\dagger$?

**回顾感知器算法：**

- 对样本的特征向量 **x** 和权向量 **w** 增广化
- 初始化权向量 $w_0$（例如： $w_0 = \mathbf{0}$）
- $for\ t = 0,1,2,\dots$ （$t$ 代表迭代次数）
  - ① 对某些样本 $n$，通过下式对权向量 $w_t$ 进行更新：

$$w_{t+1} = w_t + 1 \cdot (\llbracket \operatorname{sign}(w_t^T \mathbf{x}_{n(t)}) \neq y_n \rrbracket \mathbf{x}_{n(t)})$$

…直到满足停止条件，此时的 $w_{t+1}$ 作为学到的 $g$

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

**迭代优化：**

- $for\ t = 0,1,2,\dots$ （$t$ 代表迭代次数）

$$w_{t+1} = w_t + \eta \cdot \mathbf{v}$$

…直到满足停止条件，此时的 $w_{t+1}$ 作为学到的 $g$



by Lili Jiang

**迭代优化：**

假设 $L_{in}(w)$ 是单变量 $w$ 的函数 　　　　 $w^* = arg\min_{w} L_{in}(w)$



Loss $L_{in}(w)$

η 称为学习率 "*learning rate*"

- 随机地选择一初始值 $w_0$
- 计算 $\frac{dL_{in}}{dw}|_{w=w_0}$ 　　　 $w_1 \leftarrow w_0 - \eta \frac{dL_{in}}{dw}|_{w=w_0}$
- 计算 $\frac{dL_{in}}{dw}|_{w=w_1}$ 　　　 $w_2 \leftarrow w_1 - \eta \frac{dL_{in}}{dw}|_{w=w_1}$

...... 很多次迭代

$-\eta \frac{dL_{in}}{dw}|_{w=w_0}$

$w_0$　$w_1$　$w_2$　$w_t$　　　　　　　$w$

# 梯度下降法

## 3.3 梯度下降法

**训练(Training)**



输入图像
特征表达

回归模型 $g$

预测结果 $\hat{y}_n$

真实结果 $y_n$ $f$

损失函数

误差度量

回归模型参数空间 $\mathcal{H}(\cdot)$

参数优化 $\mathcal{A}$

**梯度下降法：**

$$\boldsymbol{w}_{t+1} \leftarrow \boldsymbol{w}_t - \eta \nabla L_{in}(\boldsymbol{w}_t)$$

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

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

$$L_{in}(\boldsymbol{w}) = \frac{1}{N}\sum_{n=1}^{N}(\boldsymbol{w}^T\mathbf{x}_n - y_n)^2$$

$$\nabla L_{in}(\boldsymbol{w}) = \frac{2}{N}\sum_{n=1}^{N}(\boldsymbol{w}^T\mathbf{x}_n - y_n)\mathbf{x}_n$$

## 梯度下降法实现线性回归

- 初始化权向量 $\boldsymbol{w}_0$

- $for\ t = 0,1,2,\dots$ ($t$ 代表迭代次数)

  ① 计算梯度: $\nabla L_{in}(\boldsymbol{w}) = \sum_{n=1}^{N}(\boldsymbol{w}^T \mathbf{x}_n - y_n)\mathbf{x}_n$

  ② 对权向量 $\boldsymbol{w}_t$ 进行更新: $\boldsymbol{w}_{t+1} \leftarrow \boldsymbol{w}_t - \eta \nabla L_{in}(\boldsymbol{w}_t)$

...直到 $\nabla L_{in}(\boldsymbol{w}) = \mathbf{0}$ ，或者迭代足够多次数

返回最终的 $\boldsymbol{w}_{t+1}$ 作为学到的 $g$

---

### 3.3 梯度下降法

**训练(Training)**

输入图像
特征表达

回归模型 $g$

预测结果 $\hat{y}_n$

真实结果 $y_n$  $f$

损失函数

误差度量

回归模型参数空间 $\mathcal{H}(\cdot)$

参数优化 $\mathcal{A}$

**梯度下降法:**

$$\nabla L_{in}(\boldsymbol{w}) = \sum_{n=1}^{N}(\boldsymbol{w}^T \mathbf{x}_n - y_n)\mathbf{x}_n$$

$$\boldsymbol{w}_{t+1} \leftarrow \boldsymbol{w}_t - \eta \nabla L_{in}(\boldsymbol{w}_t)$$

➤ *问题1: 学习率 $\eta$ 如何取值?*
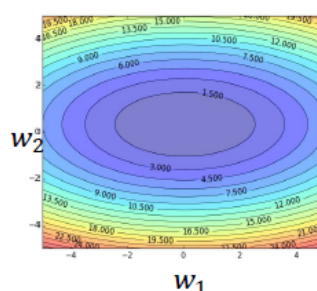
人工智能与自动化学院

---

**损失函数曲线的可视化:**



单变量 $w$ 的
损失函数曲线

两个变量 $w$ 的
损失函数曲线

损失函数曲线的等高线表示
（红色陡峭，蓝色平缓）

# AdaGrad 自适应梯度下降法>

**学习率 $\eta$ 的取值讨论：**



损失函数与迭代次数的关系曲线

梯度幅值与迭代次数的关系曲线

**此时已经到达谷底了？**



损失函数曲线的等高线表示

**学习率 $\eta$ 的取值讨论：**



梯度变化剧烈，需要小的学习率

梯度变化平缓，需要大的学习率

$$w_{t+1} \leftarrow w_t - \eta \nabla L_{in}(w_t)$$

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

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

与 $i$ 有关，
与迭代时刻的梯度有关

## 3.3 梯度下降法

**Root Mean Square：**

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

$$w_{i,1} \leftarrow w_{i,0} - \frac{\eta}{\sigma_{i,0}} \frac{\partial L_{in}}{\partial w_{i,0}} \qquad \sigma_{i,0} = \sqrt{(\frac{\partial L_{in}}{\partial w_{i,0}})^2}$$

$$w_{i,2} \leftarrow w_{i,1} - \frac{\eta}{\sigma_{i,1}} \frac{\partial L_{in}}{\partial w_{i,1}} \qquad \sigma_{i,1} = \sqrt{\frac{1}{2}[(\frac{\partial L_{in}}{\partial w_{i,0}})^2 + (\frac{\partial L_{in}}{\partial w_{i,1}})^2]}$$

$$w_{i,3} \leftarrow w_{i,2} - \frac{\eta}{\sigma_{i,2}} \frac{\partial L_{in}}{\partial w_{i,2}} \qquad \sigma_{i,2} = \sqrt{\frac{1}{3}[(\frac{\partial L_{in}}{\partial w_{i,0}})^2 + (\frac{\partial L_{in}}{\partial w_{i,1}})^2 + (\frac{\partial L_{in}}{\partial w_{i,2}})^2]}$$

$$w_{i,t+1} \leftarrow w_{i,t} - \frac{\eta}{\sigma_{i,t}} \frac{\partial L_{in}}{\partial w_{i,t}} \qquad \sigma_{i,t} = \sqrt{\frac{1}{t+1} \sum_{t=0}^{t} (\frac{\partial L_{in}}{\partial w_{i,t}})^2}$$

自适应动态学习率

（Learning rate adapts dynamically）:

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

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

**AdaGrad**

曲线平缓，$\sigma_{i,t}$小
学习率(迭代步长)大

$\frac{\partial L_{in}}{\partial w_{i,t-1}}$

$\frac{\partial L_{in}}{\partial w_{i,t}}$

$w_i$

曲线陡峭，$\sigma_{i,t}$大
学习率(迭代步长)小

$\frac{\partial L_{in}}{\partial w_{i,t-1}}$

$\frac{\partial L_{in}}{\partial w_{i,t}}$

$w_i$

自适应动态学习率(Learning rate adapts dynamically):

梯度变化平缓，
需要大的学习率

梯度变化剧烈，
需要小的学习率

$w_2$

$w_1$

**RMSProp**

**RMSProp：**

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

$$w_{i,1} \leftarrow w_{i,0} - \frac{\eta}{\sigma_{i,0}} \frac{\partial L_{in}}{\partial w_{i,0}} \qquad \sigma_{i,0} = \sqrt{(\frac{\partial L_{in}}{\partial w_{i,0}})^2}$$

$$0 < \alpha < 1$$

$$w_{i,2} \leftarrow w_{i,1} - \frac{\eta}{\sigma_{i,1}} \frac{\partial L_{in}}{\partial w_{i,1}} \qquad \sigma_{i,1} = \sqrt{\alpha(\sigma_{i,0})^2 + (1-\alpha)(\frac{\partial L_{in}}{\partial w_{i,1}})^2}$$

$$w_{i,3} \leftarrow w_{i,2} - \frac{\eta}{\sigma_{i,2}} \frac{\partial L_{in}}{\partial w_{i,2}} \qquad \sigma_{i,2} = \sqrt{\alpha(\sigma_{i,1})^2 + (1-\alpha)(\frac{\partial L_{in}}{\partial w_{i,2}})^2}$$
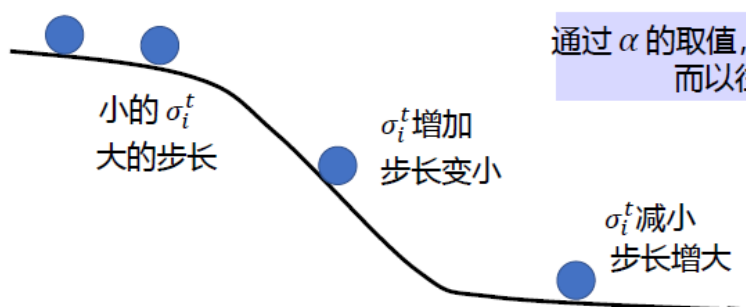
$$w_{i,t+1} \leftarrow w_{i,t} - \frac{\eta}{\sigma_{i,t}} \frac{\partial L_{in}}{\partial w_{i,t}} \qquad \sigma_{i,t} = \sqrt{\alpha(\sigma_{i,t-1})^2 + (1-\alpha)(\frac{\partial L_{in}}{\partial w_{i,t}})^2}$$

**RMSProp：**

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

$$\frac{\partial L_{in}}{\partial w_{i,0}}, \quad \frac{\partial L_{in}}{\partial w_{i,1}}, \quad ..., \quad \frac{\partial L_{in}}{\partial w_{i,t-1}},$$

$$0 < \alpha < 1$$

$$\sigma_{i,t} = \sqrt{\alpha(\sigma_{i,t-1})^2 + (1-\alpha)(\frac{\partial L_{in}}{\partial w_{i,t}})^2}$$
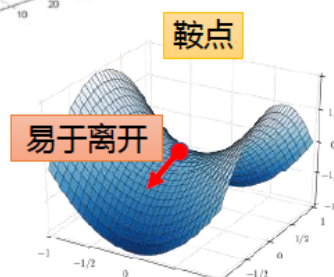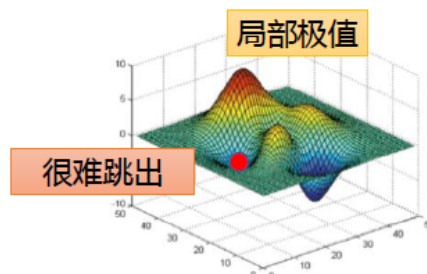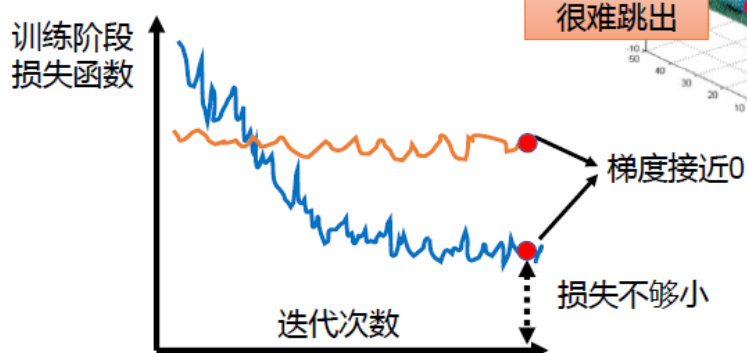
通过 $\alpha$ 的取值，使得当前的梯度影响更大，而以往的梯度影响较小



小的 $\sigma_i^t$ 大的步长

$\sigma_i^t$ 增加 步长变小

$\sigma_i^t$ 减小 步长增大

**问题2：梯度为0就能得到全局最优解吗？** ×

梯度为0时没能得到最佳解



局部极值

很难跳出

梯度接近0

鞍点

易于离开

训练阶段 损失函数

迭代次数

损失不够小

# Momentum 动量法梯度下降

## 梯度较小时几乎停止更新合理吗？



当一个球从山坡上滚下时 …

这样的物理特性能用于梯度下降吗？

单变量w

## 用向量几何示意一般的梯度下降过程



→ Gradient
→ Movement

初始值为 $w_0$

计算梯度：$\nabla L_{in}(w_0)$

权向量更新：$w_1 \leftarrow w_0 - \eta \nabla L_{in}(w_0)$

计算梯度：$\nabla L_{in}(w_1)$

权向量更新：$w_2 \leftarrow w_1 - \eta \nabla L_{in}(w_1)$

## 结合动量特性的梯度下降法(Gradient Descent + Momentum)

动量:上一步的动量减去当前梯度



→ 梯度向量
→ 动量
⋯⋯ 上一步动量

初始值为 $w_0$

动量 $m_0 = 0$

计算梯度：$\nabla L_{in}(w_0)$

动量 $m_1 = \lambda m_0 - \eta \nabla L_{in}(w_0)$

权向量更新 $w_1 \leftarrow w_0 + m_1$

计算梯度：$\nabla L_{in}(w_1)$

动量 $m_2 = \lambda m_1 - \eta \nabla L_{in}(w_1)$

权向量更新 $w_2 \leftarrow w_1 + m_2$

权向量更新时不仅考虑负梯度方向
还同时要考虑上一步的动量

## 结合动量特性的梯度下降法(*Gradient Descent + Momentum*)

动量: 当前负梯度与上一步动量之和

loss

$\longrightarrow$ Negative of $\partial L / \partial w$

$\cdots\blacktriangleright$ Last Movement

$\longrightarrow$ Real Movement

$\partial L/\partial w = 0$

# Adam 亚当优化器

## Adam: RMSProp + Momentum

**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. $g_t^2$ 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
  $m_0 \leftarrow 0$ (Initialize 1st moment vector) → for momentum
  $v_0 \leftarrow 0$ (Initialize 2nd moment vector) → 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 $t$)
    $m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$ (Update biased first moment estimate)
    $v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$ (Update biased second raw moment estimate)
    $\hat{m}_t \leftarrow m_t/(1 - \beta_1^t)$ (Compute bias-corrected first moment estimate)
    $\hat{v}_t \leftarrow v_t/(1 - \beta_2^t)$ (Compute bias-corrected second raw moment estimate)
    $\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \hat{m}_t/(\sqrt{\hat{v}_t} + \epsilon)$ (Update parameters)
  **end while**
  **return** $\theta_t$ (Resulting parameters)

$$\nabla L_{in}(\boldsymbol{w}) = \sum_{n=1}^{} (\boldsymbol{w}^T \mathbf{x}_n - y_n)\mathbf{x}_n$$

$$m_{i,t+1} = \lambda m_{i,t} - \frac{\eta}{\sigma_{i,t}} \frac{\partial L_{in}}{\partial w_{i,t}}$$

$$w_{i,t+1} \leftarrow w_{i,t} + m_{i,t+1}$$

## 问题3：训练样本批量大小的影响？

# batch_size 的影响

## 梯度下降法实现线性回归

- 初始化权向量 $\boldsymbol{w}_0$
- $\boldsymbol{for}\ t = 0,1,2,\ldots$ ($t$ 代表迭代次数)

  ① 计算梯度：$\nabla L_{in}(\boldsymbol{w}) = \dfrac{2}{N}\sum_{n=1}^{N}(\boldsymbol{w}^T\mathbf{x}_n - y_n)\mathbf{x}_n$  每次迭代N个样本均要计算，时间复杂度与Pocket算法相似

  ② 对权向量 $\boldsymbol{w}_t$ 进行更新：$\boldsymbol{w}_{t+1} \leftarrow \boldsymbol{w}_t - \eta\nabla L_{in}(\boldsymbol{w}_t)$

...直到 $\nabla L_{in}(\boldsymbol{w}) = \boldsymbol{0}$，或者迭代足够多次数

返回最终的 $\boldsymbol{w}_{t+1}$ 作为学到的 $g$

## 梯度下降法实现线性回归

- 初始化权向量 $\boldsymbol{w}_0$
- $\boldsymbol{for}\ t = 0,1,2,\ldots$ ($t$ 代表迭代次数)

  ① 计算梯度：$\nabla L_{in}(\boldsymbol{w}) = (\boldsymbol{w}^T\mathbf{x}_n - y_n)\mathbf{x}_n$  随机梯度下降法 (Stochastic Gradient Descent) (SGD)

  ② 对权向量 $\boldsymbol{w}_t$ 进行更新：$\boldsymbol{w}_{t+1} \leftarrow \boldsymbol{w}_t - \eta\nabla L_{in}(\boldsymbol{w}_t)$

...直到 $\nabla L_{in}(\boldsymbol{w}) = \boldsymbol{0}$，或者迭代足够多次数

返回最终的 $\boldsymbol{w}_{t+1}$ 作为学到的 $g$

## 迭代过程中一次计算样本多少(batch)对梯度下降的影响

假设一共有20个样本



Batch size = N (Full batch)

计算20个样本才更新

20个样本均参与计算后才更新

一个样本更新一次

Batch size = 1

20次更新以后

计算每个样本后就更新

# 迭代过程中一次计算样本多少(batch)对梯度下降的影响

每次更新时花费的时间

MNIST: digit
classification

Having limitation

Parallel computing

Tesla V100 GPU

批量(batch size)大在计算梯度时并不意味着会花费更多的时间，除非batch size 太大

# 批量(batch)用于梯度下降

$$\boldsymbol{w}^* = arg\min_{\boldsymbol{w}} L$$

➤ 随机设置一个初始向量 $\boldsymbol{w}_0$

➤ 计算梯度 $\nabla L_1(\boldsymbol{w}_0)$，更新： $\boldsymbol{w}_1 \leftarrow \boldsymbol{w}_0 - \eta \nabla L_1(\boldsymbol{w}_0)$    $L_1$

➤ 计算梯度 $\nabla L_2(\boldsymbol{w}_1)$，更新： $\boldsymbol{w}_2 \leftarrow \boldsymbol{w}_1 - \eta \nabla L_2(\boldsymbol{w}_1)$    $L_2$

➤ 计算梯度 $\nabla L_3(\boldsymbol{w}_2)$，更新： $\boldsymbol{w}_3 \leftarrow \boldsymbol{w}_2 - \eta \nabla L_3(\boldsymbol{w}_2)$    $L_3$

B { batch    L

batch

batch    N

batch

**1 epoch** = 所有的batch遍历一遍 ➡ 每一次epoch后进行Shuffle

# 批量(batch)大小对训练过程速度的影响

一次更新所花费的时间      一个epoch所花费的时间

60 updates      slower

     faster

60000 updates in one epoch

较小的批量(batch)做完一次epoch时需要花费更多的时间

# 批量(batch)大小对分类性能的影响

## MNIST

Accuracy vs. Batch Size



## CIFAR-10

Accuracy vs. Batch Size



较小的批量获得更好的性能

为什么批量大了性能会下降呢?

# 批量(batch)大小对分类性能的影响

SB = 256

LB = 0.1 x data set

| Name | Network Type | Data set |
|------|--------------|----------|
| $F_1$ | Fully Connected | MNIST (LeCun et al., 1998a) |
| $F_2$ | Fully Connected | TIMIT (Garofolo et al., 1993) |
| $C_1$ | (Shallow) Convolutional | CIFAR-10 (Krizhevsky & Hinton, 2009) |
| $C_2$ | (Deep) Convolutional | CIFAR-10 |
| $C_3$ | (Shallow) Convolutional | CIFAR-100 (Krizhevsky & Hinton, 2009) |
| $C_4$ | (Deep) Convolutional | CIFAR-100 |

On Large-Batch Training for Deep Learning: Generalization Gap and Sharp Minima

https://arxiv.org/abs/1609.04836

| Name | Training Accuracy | | Testing Accuracy | |
|------|------|------|------|------|
| | SB | LB | SB | LB |
| $F_1$ | $99.66\% \pm 0.05\%$ | $99.92\% \pm 0.01\%$ | $98.03\% \pm 0.07\%$ | $97.81\% \pm 0.07\%$ |
| $F_2$ | $99.99\% \pm 0.03\%$ | $98.35\% \pm 2.08\%$ | $64.02\% \pm 0.2\%$ | $59.45\% \pm 1.05\%$ |
| $C_1$ | $99.89\% \pm 0.02\%$ | $99.66\% \pm 0.2\%$ | $80.04\% \pm 0.12\%$ | $77.26\% \pm 0.42\%$ |
| $C_2$ | $99.99\% \pm 0.04\%$ | $99.99\% \pm 0.01\%$ | $89.24\% \pm 0.12\%$ | $87.26\% \pm 0.07\%$ |
| $C_3$ | $99.56\% \pm 0.44\%$ | $99.88\% \pm 0.30\%$ | $49.58\% \pm 0.39\%$ | $46.45\% \pm 0.43\%$ |
| $C_4$ | $99.10\% \pm 1.23\%$ | $99.57\% \pm 1.84\%$ | $63.08\% \pm 0.5\%$ | $57.81\% \pm 0.17\%$ |

较小的批量在测试数据集上也能获得更好的性能

# 批量(batch)大小对梯度下降优化时的影响总结

| | 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 WIN |
| 梯度的特点 | Noisy | Stable |
| 优化性能 | Better WIN | Worse |
| 泛化性能 | Better WIN | Worse |

批量大小(batch size)作为超参数(hyperparameter)由算法设计者确定

**训练(Training)**

| 输入图像 特征表达 |
|---|

回归模型 $g$

预测结果 $\hat{y}_n$

真实结果 $y_n$ $f$

损失函数

误差度量

回归模型参数空间 $\mathcal{H}(\cdot)$

$\mathcal{A}$ 参数优化

**随机梯度下降法：**

$$\nabla L_{in}(w) = \sum_{n=1}^{B}(w^T \mathbf{x}_n - y_n)\mathbf{x}_n$$

$$\boldsymbol{m}_{i,t+1} = \lambda \boldsymbol{m}_{i,t} - \frac{\eta}{\sigma_{i,t}}\frac{\partial L_{in}}{\partial w_{i,t}}$$

$$\boldsymbol{w}_{i,t+1} \leftarrow \boldsymbol{w}_{i,t} + \boldsymbol{m}_{i,t+1}$$

- 问题1：学习率 $\eta$ 如何取值？
- 问题2：梯度为0就能得到最佳解？
- 问题3：训练样本批量大小的影响？

# 小结

3.1 线性回归问题

*模型的输出为实数值，有众多应用场景*

3.2 线性回归算法

*损失函数为均方误差时，可通过求解广义逆得到解析解*

3.3 梯度下降法

*迭代优化，更一般的损失函数；固定学习率、AdaGrad、RMSProp、动量(Momentum)、Adam、SGD、批量大小(batch size)*

# 作业

# 手写作业

2，根据向量或矩阵的计算性质，证明：

$$\|\mathbf{X}w - Y\|^2 = w^T\mathbf{X}^T\mathbf{X}w - 2w^T\mathbf{X}^TY + Y^TY$$

$T_2$. 解: 证明: $\|Xw-Y\|^2 = w^T X^T X w - 2w^T X^T Y + Y^T Y$

$$\|Xw-Y\|^2 = (Xw-Y)^T (Xw-Y)$$

$$= (w^T X^T - Y^T)(Xw-Y)$$

$$= w^T X^T X w - w^T X^T Y - \boxed{Y^T X w} + Y^T Y$$

$$= w^T X^T X w - w^T X^T Y - \boxed{(Xw)^T Y} + Y^T Y$$

$$= w^T X^T X w - 2w^T X^T Y + Y^T Y.$$