# logistic regression.

- Sigmoid function: $g(z) = \dfrac{1}{1+e^{-z}}$

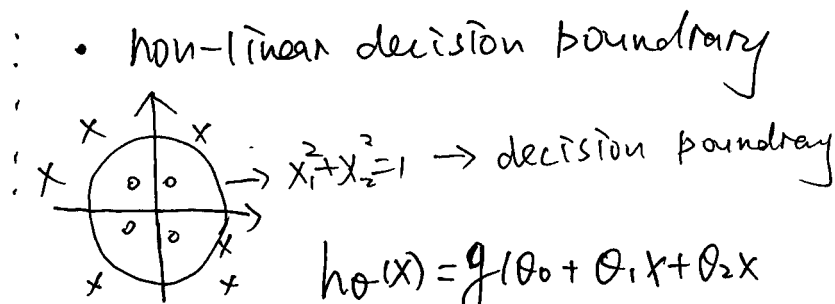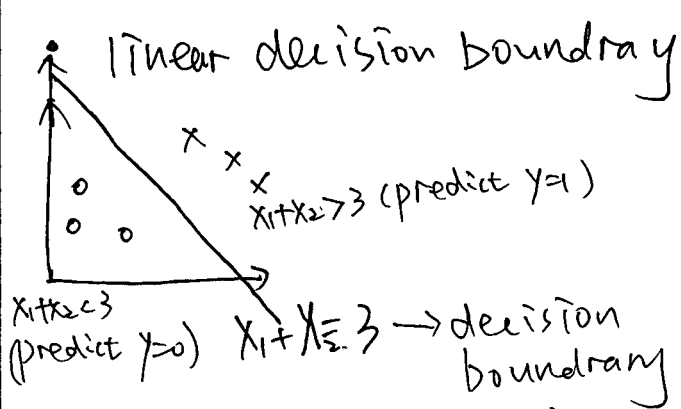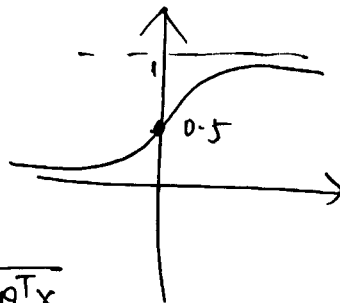$$\Updownarrow$$

$$h_\theta(X) = g(\theta^T X) = \dfrac{1}{1+e^{-\theta^T X}}$$

$$= \text{estimated prob for } y=1, \text{ based on input } X$$

eg. $X = \begin{pmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 1 \\ sex \\ height \\ weight \end{pmatrix}$  若 $h_\theta(X) = 0.7 \Leftrightarrow P(y=1|X,\theta) = 0.7$

这个人有 70% chance 有 8号身材

- linear decision boundray :



$x_1 + x_2 > 3$ (predict $y=1$)

$x_1 + x_2 < 3$ (predict $y=0$)  $x_1 + x_2 = 3 \to$ decision boundray

eg. $h_\theta(X) = g(\theta_0 + \theta_1 X + \theta_2 X)$

$\theta_0 = -3 \quad \theta_1 = 1 \quad \theta_2 = 1$

即 $h_\theta(X) = X_1 + X_2 - 3$

predict $y=1$ if $-3 + X_1 + X_2 \geq 0$

- non-linear decision boundray



$X_1^2 + X_2^2 = 1 \to$ decision boundray

$$h_\theta(X) = g(\theta_0 + \theta_1 X + \theta_2 X + \theta_3 X_1^2 + \theta_4 X_2^2)$$

$$= g(X_1^2 + X_2^2)$$

predict $y=1$ ~~if $-1+X_1^2+X_2^2$~~

if $-1 + X_1^2 + X_2^2 \geq 0$

- core: how to fit $\theta = [\theta_0 \cdots \theta_n]$ ?   input: $\{(X_1,y_1) \cdots (X_m,y_m)\}$   $X \in R^{n+1}$

不能用 最小二乘. ∵ $J(\theta) = \dfrac{1}{m} \sum_1^m \dfrac{1}{2}(h_\theta(X_i) - y_i)^2$

∵ $J(\theta)$ 会是 non-convex. 

- $\text{cost}(h_\theta(X), y) = \begin{cases} -\log(h_\theta(X)) & y=1 \\ -\log(1-h_\theta(X)) & y=0 \end{cases}$

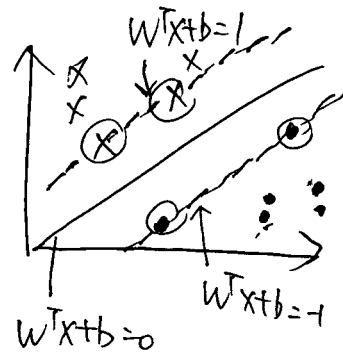① 若 label = 1, 而 $h_\theta(X) = 0$,  ← $y=1$的 prob

则会给一个很大的 penalty

② 若 label = 0, 而 $h_\theta(X) = 1$

则会给一个很大 penalty

- Why choose this ?

再加入 正则项得 损失函数 $J(\theta) = -\frac{1}{m}\left[\sum_{i=1}^{m} y_i \log(h_\theta(x_i)) + (1-y_i)\log(1-h_\theta(x_i))\right]$
$+\frac{\lambda}{2m}\sum_{j=1}^{n}\theta_j^2$

$\frac{\lambda}{2m}\sum_{j=1}^{n}\theta_j^2$     用 gb fit $\theta$. s.t. $\min_\theta J(\theta)$

repeat: $\theta_j = \theta_j - \alpha \cdot \frac{dJ(\theta)}{d\theta_j}$
↑
步长.

$\boxed{\text{SVM}}$ 给定 $D = \{(x_1,y_1)\cdots(x_m,y_m)\}$, $y_i \in \{-1,1\}$, 找到超平面.
将不同类别分开. 正确的分类条件 $\begin{cases} w^Tx_i+b>0 & y_i=1 \\ w^Tx_i+b<0 & y_i=-1 \end{cases}$ $(w^Tx+b=0)$

SVM 长3个心眼. 条件变为 $\begin{cases} w^Tx_i+b\geq1 & y_i=1 \\ w^Tx_i+b\leq-1 & y_i=-1 \end{cases}$

间隔带宽 $d = \frac{2}{\|w\|}$

So. 求 $\max\limits_{w,b}\frac{2}{\|w\|}$, s.t. $y_i(w^Tx_i+b)\geq1$ $i=1\cdots m$

也就是 $\min \frac{\|w\|^2}{2}$, s.t. $y_i(w^Tx_i+b)\geq1$ $i=1\cdots m$. —primal

∵无法直接解 $w,b$
用对偶问题求解 $w,b$, 引入 拉格朗日函数 $L(w,b,\alpha) = \frac{1}{2}\|w\|^2 + \sum_{i=1}^{m}\alpha_i(1-y_i(w^Tx_i+b))$

$\frac{\partial L(w,b,\alpha)}{\partial w}=0 \Rightarrow w = \sum_{i=1}^{m}\alpha_i y_i x_i$

$\frac{\partial L(w,b,\alpha)}{\partial b}=0 \Rightarrow b = \sum_{i=1}^{m}\alpha_i y_i$

$\left.\right\}$代入 $L(w,b,\alpha)$ 得 dual 问题

dual: $\max\limits_{\alpha} \sum_{i=1}^{m}\alpha_i - \frac{1}{2}\sum_{i=1}^{m}\sum_{j=1}^{m}\alpha_i\alpha_j y_i y_j x_i^T x_j$  s.t. $\sum_{i=1}^{m}\alpha_i y_i=0$. $\alpha_i\geq0$ $i=1\cdots m$

objective: $\arg\max\limits_{\alpha}(\sum_i\alpha_i - \frac{1}{2}\sum_i\sum_j\alpha_i\alpha_j y_i y_j x_i^T x_j)$ s.t. $\sum_i\alpha_i y_i=0$
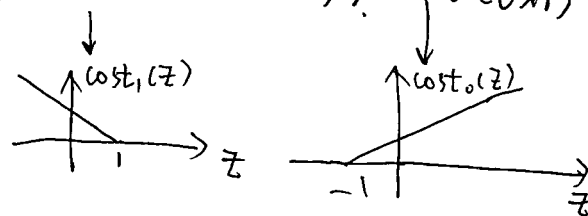
求出 $\alpha^*$后
$\Rightarrow w^* = \sum_i\alpha_i^* y_i x_i$

❈ SVM 的 lost function: $\min C\sum_{i=1}^{m}[y_i \text{Cost}_1(\theta^Tx_i)+(1-y_i)\text{Cost}_0(\theta^Tx_i)]$
hinge loss: 分类对 损失为0 $+\frac{1}{2}\sum_{j=1}^{m}\theta_j^2$

C: 软间隔. 若 希望 分错点 越少越好, 则大 C.
若C 很小, 对 outlier 不理会.
∴ SVM 对 outlier 不敏感. LR 敏感

梯度下降: 以 $h_\theta = \sum_{j=0}^{n} \theta_j x_j$ 为例.

损失函数 $J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2$

① BGD: 批量梯度下降.    目标: 求解 weights 使 误差函数 $J$ 于最小

先随机取 weights. 不断 更新 weights s.t. $J$ 减小.

$$\theta_j = \theta_j - \alpha \cdot \frac{\partial J(\theta)}{\partial \theta_j}$$

↑ learning rate: 每次向着 $J$ 最陡峭的方向迈步

对点 $(x, y)$

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{\partial (\frac{1}{2}(h_\theta(x) - y)^2)}{\partial \theta_j} = 2 \cdot \frac{1}{2}(h_\theta(x) - y) \cdot \frac{\partial}{\partial \theta_j}(h_\theta(x) - y)$$

$$= (h_\theta(x) - y) \cdot \frac{\partial (\sum_{i=0}^{n} \theta_i x_i - y)}{\partial \theta_j} = (h_\theta(x) - y) x_j$$

则 对 all data points: 偏导累和 $\frac{\partial J(\theta)}{\partial \theta_j} = -\frac{1}{m} \sum_{i=1}^{m} (y^i - h_\theta(x^i)) x_j^i$

repeat {
$\theta_j' = \theta_j + \frac{1}{m} \sum_{i=1}^{m} (y^i - h_\theta(x^i)) x_j^i$

for every $j = 0 \dots n$     #∵ 有 n 个参数
}

• 每次训练都用到全部 data 很耗时.
• 但 迭代次数较少

② SGD: 随机梯度下降
(为了解决 BGD 太慢)
for $i = 1, \dots, m$
$\theta_j' = \theta_j + (y^i - h_\theta(x^i)) x_j^i$ (for $j = 0 \dots n$)

sgd: 通过每个样本来更新迭代一次. 但 噪音比 bgd 多, s.t. sgd
不是每次都朝着 整体一致的方向.  次数比 bgd 多.

缺点: 更新方向完全依赖 方向 batch, 更新不稳定.  ∴ 可以引入 momentum (更新时一定 程度沿用之前更新的方向)
( BGD 和 SGD 不折衷 更重要的! )
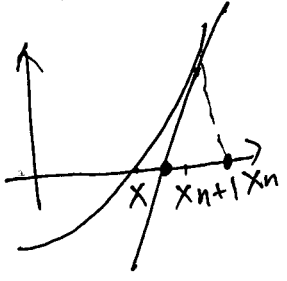
③ min_batch gd = MBGD

s.t. 快且 high acc.

Repeat { for $i = 1, 11, 21 \dots 991$ {
$\theta_j' = \theta_j - \alpha \cdot \frac{1}{10} \sum_{k=i}^{i+9} [h_\theta(x^{(k)}) - y^{(k)}] x_j^{(k)}$

for every $j = 0 \dots n$)
}
}

缺点: 不能保证 很好的收敛性.

牛顿法. 应用 { 求方程根
                   求最优化办法



1) 迭代 求根.    $f(x)=0$

$f(x_0+\Delta x) = f(x_0) + f'(x_0)*\Delta x$    即 $f(x) = f(x_0) + f'(x_0)(x-x_0)$
                                                        一阶 Talyor

令 $f(x)=0$    得迭代公式: $x = x_0 - f(x_0)/f'(x_0)$

$$\Downarrow$$

$$x_{n+1} = x_n - f(x_n)/f'(x_n)$$    逐步迭代 求最优解

2) 推广到 非线性最优化问题

2阶Taloyer:  $f(x) = f(x_0) + f'(x_0)(x-x_0) + \frac{1}{2}f''(x_0)(x-x_0)^2 + \underline{O((x-x_0)^3)}$    ~~$O((x-x_0)^3)$~~    ignore.

代入 $x = x_0 + \Delta x$    得  $f(x) = f(x_0+\Delta x) = f(x_0) + f'(x_0)\cdot\Delta x + \frac{1}{2}f''(x_0)(\Delta x)^2$

$f'(x) = f'(x_0+\Delta x) = 0 \Rightarrow [f(x_0) + f'(x_0)(x-x_0) + \frac{1}{2}f''(x_0)(x-x_0)^2]' = 0$

$\Rightarrow f'(x_0) + f''(x_0)(x-x_0) + \frac{1}{2}f''(x_0)\cdot 2(x-x_0) = 0$

$\Rightarrow f'(x_0) + f''(x_0)(x-x_0) = 0$    $\Rightarrow x = x_0 - f'(x_0)/f''(x_0)$    $\Rightarrow x_{n+1} = x_n - \frac{f'(x_n)}{f''(x_n)}$

用了二阶导, 收敛更快

[Hessian]  对于多变量问题. 牛顿法演变为  $x_{n+1} = x_n - \frac{J_f(x_n)}{H(x_n)}$

$$= x_n - H^{-1}(x_n)\cdot J_f(x_n)$$

$J$: 雅克比矩阵. $J_f(x_n) = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \vdots & & \\ \frac{\partial y_m}{\partial x_1} & \cdots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}$

H: Hessian矩阵.
~~牛顿~~ 起到了控制步长的    $H(f) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 x_n} \\ \vdots & & & \\ \frac{\partial^2 f}{\partial x_n x_1} & & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$

作用、特征值 更新步长

$H = E[\lambda_1 \cdots \lambda_n]E^T$    E是特征
                                           所对应的特征矩阵

H 非正定时. 无法收敛

H 矩阵维度过大. 带来巨大计算量    $x=0$时:
例如 一元函数 $f(x)=x^2$. $f'(x)=2x$. $f''(x)=2$    $\because f'(x)=0$.  而 $f''(x)>0$. $\therefore x=0$ 为极小值
类似的. $|H|=0$ 时. 无法确定函数是否取得极值.    $\therefore$ { 正定: 临界点 极小 $(\lambda_i \, 全>0)$
                                                                                          负定: 临界点 极大 $(\lambda_i<0)$
                                                                                          不定: 鞍点 $(o.w.)$

Momentum. 通过加入 $\gamma V_{t-1}$、加速 sgd、且抑制震荡

$V_t = \gamma V_{t-1} + \eta \nabla_\theta J(\theta)$, $\theta = \theta - V_t$ 例如 ⊙ 从山顶滚下来，有阻力 滚的慢

加入的这一项可以使 梯度方向不变的维度\(变快，梯度方向

改变的维度更新速度变慢。这样可以 加速、并 ↓震荡.      $\gamma = 0.9$

Adagrad: 对低频多数做较大更新.高频的较小更新. ∴对稀疏数据表现好

Adadelta: 对 Adagrad 改进

$$\boxed{手推公式}$$

xgboost

· 如何集成:

  基础: decision tree

  集成: 一个个加的

  $\hat{y_i}^{(0)} = 0$

  $\hat{y_i}^{(1)} = f_1(x_i) = \hat{y_i}^{(0)} + f_1(x_i)$

  $\hat{y_i}^{(2)} = f_1(x_i) + f_2(x_i) = \hat{y_i}^{(1)} + f_2(x_i)$

  $\hat{y_i}^{(t)} = \sum\limits_{k=1}^{t} f_k(x_i) = \underline{\hat{y_i}^{(t-1)}} + f_t(x_i)$ 加入一个新函数预测

  ↑        保留前 t-1 轮的模型预测

  前轮的预测模型

· 每一轮加入的模型如何构造？ 方案:优化了我的 objective

  $Obj^{(t)} = \sum\limits_{i=1}^{n} L(y_i, \hat{y_i}^{(t)}) + \sum\limits_{i=1}^{t} \Omega f(i)$      $\Omega(f_t) = \gamma T + \frac{1}{2}\lambda \sum\limits_{j=1}^{T} w_j^2$

  $= \sum\limits_{i=1}^{n} L(y_i, \hat{y_i}^{(t-1)} + f_t(x_i)) + \underbrace{\sum \Omega f_t}_{} \Omega(f_t) + C$        $\underbrace{叶子个数}_{限制叶子个数}$  $\underbrace{}_{权重正则化}$

· 目标函数:                                                        防止权过大

  ① 保证 bias 小   ② 保证树模型最精简

Talyor 展开: $f(x + \delta x) \approx f(x) + f'(x)\triangle x + \frac{1}{2}f''(x)\delta x^2$

  $g_i =$ 一阶导  $h_i =$ 二阶导 $= \partial^2_{\hat{y}^{(t-1)}} L(y_i, \hat{y}^{(t-1)})$       新增

  $= \partial_{\hat{y}^{(t-1)}} L(y_i, \hat{y}^{(t-1)})$

  $Obj^{(t)} \approx \sum\limits_{i=1}^{n} [L(y_i, \hat{y_i}^{(t-1)}) + g_i \cdot f_t(x_i) + \frac{1}{2}h_i \cdot f_t^2(x_i)]$

  $+ \Omega(f_t) + C$       $= \sum\limits_{i}^{n} (g_i \cdot f_t(x_i) + \frac{1}{2}h_i f_t^2(x_i)$

  $+ \Omega(f_t)$

$$Obj^{(t)} \approx \sum_{i=1}^{n} [g_i \cdot f_t(x_i) + \tfrac{1}{2} h_i f_t^2(x_i)] + \Omega(f_t) \qquad \text{n} \to \text{样本上遍历}$$

$$= \sum_{i=1}^{n} [g_i \cdot w_q(x_i) + \tfrac{1}{2} h_i w_q^2(x_i)] + \gamma T + \lambda \cdot \tfrac{1}{2} \sum_{j=1}^{T} w_j^2$$

统计每个叶子节点有多少样本

$$= \sum_{j=1}^{T} \left( \boxed{\sum_{i \in I_j} g_i} \right) w_j + \tfrac{1}{2} \left( \sum_{i \in I_j} h_i + \lambda \right) w_j^2 \right] + \gamma T$$

叶子结点上遍历. ∴ 我们只关心 最后结果

$$G_j = \sum_{i \in I_j} g_i$$
$$H_j = \sum_{i \in I_j} h_i$$

$$Obj^{(t)} = \sum_{j=1}^{T} \left( G_j \cdot w_j + \tfrac{1}{2}(H_j + \lambda) w_j^2 \right) + \gamma T$$

求导. $\dfrac{\partial \Omega(f_t)}{\partial w_j} = G_j + (H_j + \lambda) w_j$

$$= 0$$

$$\therefore Obj = -\tfrac{1}{2} \sum_{j=1}^{T} \frac{G_j^2}{H_j + \lambda} + \gamma T \quad \text{☆} \qquad \Rightarrow w_j = -\frac{G_j}{H_j + \lambda}$$

代入

structure score    分数越小. 效果越好    model performance

$$Gain = \tfrac{1}{2} \left[ \frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{(G_L + G_R)^2}{(H_L + H_R + \lambda)} \right] - \gamma \quad \text{加入新叶子节点}$$
引入的复杂度代价

↑ 左子树分数    ↑ 右子树分数    ↑ 不分割我们可以拿到的分数

---

**GBDT** : 前一轮 学习器为 $f_{t-1}(x)$, 对应的 损失函数为 $L(y, f_{t-1}(x))$

新一轮迭代: 找到一个弱分类器 $h_t(x)$, s.t. $L(y, f_{t-1}(x) + h_t(x))$ 达到最小

用 loss function 的负梯度在当前 模型的值. $-\left[ \dfrac{\partial L(y, f(x_i))}{\partial f(x_i)} \right]_{f(x) = x}$

作为回归问题中提升树算法残差的近似值.

① 初始化弱分类器, 估计使 Loss function 极小化的一个常数值, 此时树只有一个根节点. $f_0(x) = \arg\min_c \sum_{i=1}^{N} L(y_i, c)$

迭代轮数 1, 2, ... M. 计算 loss function 的负梯度值在当前模型的值, $\gamma_{mi} = -\dfrac{\partial L(y, f(x_i))}{\partial f(x_i)} \Big|_{f(x) =}$

对 $i = 1 \dots N$,

② $\begin{cases} \text{对 } \gamma_{mi} \text{拟合一个回归树, 得到第 m 棵树 叶结点, 区域 } R_{mj}. \quad j = 1 \dots J \\ \text{对 } j = 1 \dots J \quad \text{计算 } C_{mj} = \arg\min_c \sum_{x_i \in R_{mj}} L(y_i, f_{m-1}(x_i) + c) \end{cases}$

更新回归树 $f_m(x) = f_{m-1}(x) + \sum_{j=1}^{J} C_{mj} I(x \in R_{mj})$

③ 输出模型: $\hat{f}(x) = f_M(x) = \sum_{m=1}^{M} \sum_{j=1}^{J} C_{mj} I(x \in R_{mj})$