

C1-3

Machine Learning by Andrew Ng, Stanford Engineering

Xiaojie Zhou

szxjzhou@163.com

2016.8.8

第四集：牛顿方法

- 牛顿方法(Newton's Method)
- 指数分布族(Exponential Family)
- 广义线性模型(Generalized Linear Models(GLMs))

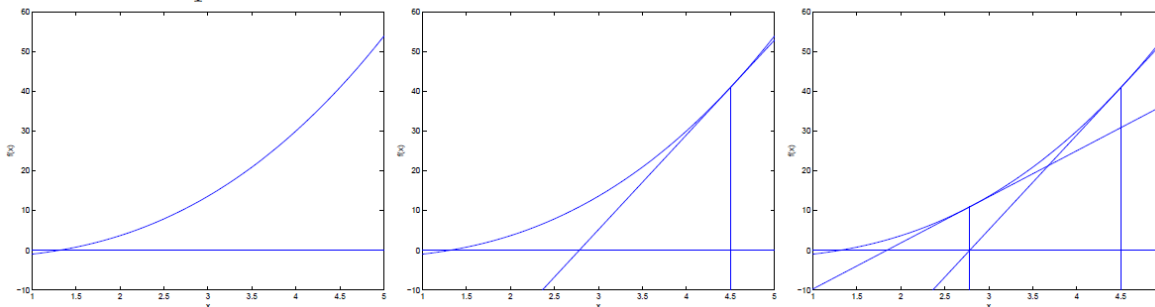
第四集：牛顿方法

- 牛顿方法

- 之前我们讲到用梯度上升法解决Logistic回归的最优化问题（极大似然），下面将讲述另外一种方法

- 对于某个实数域上的函数 f ，我们需要求 $f(\theta)=0$ 时对应的 θ 值，可以通过牛顿方法进行求解
 - 在牛顿方法中首先随机初始化 θ 的值 θ_0 ，此时我们会得到对应的函数值 $f(\theta_0)$
 - 接下来我们在 $f(\theta_0)$ 这一点处做切线，切线会与 $y=0$ 交于某一点 θ_1 ，此时又会得到新的函数值 $f(\theta_1)$ 。再在 $f(\theta_1)$ 这一点处做切线...
 - 如此迭代下去直至收敛，即可得到 $f(\theta)=0$ 时对应的 θ 值

Here's a picture of the Newton's method in action:



$$\theta := \theta - \frac{f(\theta)}{f'(\theta)}$$

第四集：牛顿方法

- 牛顿方法

- 牛顿方法给了我们一个求 $f(\theta)=0$ 时对应的 θ 值的工具

- 回到我们的问题，我们需要求 $\ell(\theta)$ 的极大值（ $\ell(\theta)$ 作为一个凹函数，其极大值等于最大值），也就相当于求 $\ell(\theta)$ 的导数为0时对应的 θ 值，由此我们可以通过下面这个更新方法迭代得到 θ 的值

$$\theta := \theta - \frac{\ell'(\theta)}{\ell''(\theta)}$$

- 牛顿方法是一种收敛速度很快的方法（二次收敛，因为考虑了二阶导数即梯度的梯度，因此通常下一次迭代的误差率接近于上一次迭代误差率的平方），对于Logistic回归、广义线性模型(GLM)作用很大
 - 上面所讲述的更新办法针对 θ 为一维参数的情况，其实对于一般意义上的牛顿方法(Newton-Raphson Method)来说更新办法如下（该更新办法中的 θ 是一个向量）

$$\theta := \theta - H^{-1} \nabla_{\theta} \ell(\theta), \quad \text{其中 } \nabla_{\theta} \ell(\theta) \text{ 为 } \ell(\theta) \text{ 的梯度, } H \text{ 为 Hessian 矩阵} \quad H_{ij} = \frac{\partial^2 \ell(\theta)}{\partial \theta_i \partial \theta_j}$$

第四集：牛顿方法

- 牛顿方法

- 比较牛顿方法和梯度下降法可以得出下面的结论

- 两者在形式上相似，其中Hessian矩阵的逆就好比梯度下降法的学习率参数 α
 - 牛顿法收敛速度相比梯度下降法更快（梯度下降法是一阶收敛，牛顿方法是二阶收敛）。这是因为牛顿方法中考虑了二阶导数即梯度的梯度，而梯度下降法只考虑了梯度，因此通常下一次迭代的误差率接近于上一次迭代误差率的平方
 - 由于海森矩阵的逆在迭代中不断减小，起到自动控制逐渐缩小步长的效果
 - 牛顿法存在的缺点是每次迭代都要重新计算一次Hessian矩阵的逆，而计算Hessian矩阵的逆是比较困难的。因此不宜用于特征很多的建模优化（因为在这些问题中 θ 值的维数很大）

第四集：牛顿方法

- 指数分布族

- 之前讲到了以线性回归的方法和Logistic回归的方法，这两种方法都涉及到对于y关于x的概率密度函数的定义方法

- 在线性回归中我们假设的是 $y = h(x) + \varepsilon$ ，其中 $h(x)$ 为 $\theta^T x$ ， ε 符合正态分布，最终得到了最小二乘衡量的结果；而在Logistic回归中我们假设的是 $y = h(x)$ ，符合二项分布，最终得到了Logistic回归的目标函数

- 实际上这两个算法都是广义线性模型这一类算法的特例

- 上面所说到的正态分布、二项分布这些概率分布都同属于指数分布族，对于指数分布族来说具有下面的形式

$$p(y; \eta) = b(y) \exp(\eta^T T(y) - a(\eta))$$

其中y为自变量（也就是之前我们常用的x）， η 称为该分布的自然参数(natural parameter/canonical parameter)， $T(y)$ 被称为充分统计量(sufficient statistic)一般情况下等于y，由于 η 为参数，因此 $\eta^T T(y)$ 一般也为实数。当我们固定下函数a,b,T后，我们可以将这个表达式看成以 η 为参数，y为自变量的一类概率分布。这里面改变了 η 以后就会得到一组新的概率分布。

第四集：牛顿方法

- 指数分布族

- 现在来看通过指数分布族的通项公式如何得到正态分布、二项分布这些概率分布

- 先来看二项分布，在二项分布中参数为 ϕ ，因变量 y 的取值要么为0要么为1

- $p(y = 1; \phi) = \phi; p(y = 0; \phi) = 1 - \phi \rightarrow p(y; \phi) = \phi^y(1 - \phi)^{1-y}$

$$\begin{aligned} p(y; \phi) &= \phi^y(1 - \phi)^{1-y} \\ &= \exp(y \log \phi + (1 - y) \log(1 - \phi)) \\ &= \exp\left(\left(\log\left(\frac{\phi}{1 - \phi}\right)\right)y + \log(1 - \phi)\right) \end{aligned}$$

第一个式子到第二个式子不是那么显然，其实对于两个式子同时取log很容易发现是相等的

- 对比指数分布族的通项公式可得 $\phi = 1/(1 + e^{-\eta})$ ，这也就是之前的Sigmoid函数，这也就印证了之前选择Sigmoid函数作为Logistic回归中 $h(x)$ 定义方法的合理性。同时我们还可得到：

$$\begin{aligned} T(y) &= y \\ a(\eta) &= -\log(1 - \phi) \\ &= \log(1 + e^{\eta}) \\ b(y) &= 1 \end{aligned}$$

$$\begin{aligned} \eta = \log \frac{\phi}{1 - \phi} &\Leftrightarrow -\eta = \log \frac{1 - \phi}{\phi} = \log\left(\frac{1}{\phi} - 1\right) \\ &\Leftrightarrow e^{-\eta} = \frac{1}{\phi} - 1 \Leftrightarrow \phi = \frac{1}{e^{-\eta} + 1} \end{aligned}$$

第四集：牛顿方法

- 指数分布族

- 接下来我们来看一下高斯分布如何表达成类似的形式

- 由于我们在线性回归里面得到最小二乘估计的过程中发现方差的取值与原优化问题无关，因此在这里我们也可以简单地将方差看成某个常数项（因此在下面我们将方差简单地设置成1，得到了如下的关系式）。此时发现这也是指数分布族的特例

$$\begin{aligned} p(y; \mu) &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(y - \mu)^2\right) \\ &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}y^2\right) \cdot \exp\left(\mu y - \frac{1}{2}\mu^2\right) \end{aligned}$$

$$\begin{aligned} \eta &= \mu \\ T(y) &= y \\ a(\eta) &= \mu^2/2 \\ &= \eta^2/2 \\ b(y) &= (1/\sqrt{2\pi}) \exp(-y^2/2) \end{aligned}$$

- 实际上就算加上方差 σ 也可以对应到指数分布族。由于涉及到两个参数，因此需要更加广义的指数分布族的定义 $p(y; \eta, \tau) = b(a, \tau) \text{Exp}[(\eta^T T(y) - a(\eta)) / c(\tau)]$ ，其中 τ 称为 dispersion parameter，在这里对应到 σ^2

第四集：牛顿方法

- 指数分布族

- 其实还有很多常见的分布都可以对应到指数分布族

- 多项式(Multinomial)分布：二项分布好比扔硬币要么为0要么为1，多项式分布相当于扔骰子，取值范围虽然也是有限数量元素集合，但不止于 $\{0,1\}$
 - 泊松(Poisson)分布：主要用于描述单位时间、面积、体积等单位内事件个数的相对频率（比如说一年内中彩票的概率、一定时间内顾客的数量）
 - Gamma分布、指数(Exponential)分布：主要用于对连续非负随机变量进行建模，和泊松分布有一定类似性（比如说泊松过程的事件间隔时间服从指数分布，从头开始到第 n 次事件发生的间隔时间服从Gamma分布），只不过这个描述的是连续的变量。比如说在公交车站上等到下一辆车的时间的概率。
 - Beta分布、狄利克雷(Dirichlet)分布：主要用于对事件成功的概率进行建模（即对概率分布进行建模）。对于硬币或者骰子这样的简单实验，我们事先能很准确地掌握系统成功的概率。然而通常情况下，系统成功的概率是未知的。比如说特殊的硬币，我们需要知道这枚硬币正面朝上的概率。然后抛10次硬币，出现5次正面，于是我们认为硬币出现正面的概率 p 最可能是0.5。但是即使硬币出现正面的概率 p 为0.4，也会出现抛10次出现5次正面的情况。因此我们并不能完全确定硬币出现正面的概率就是0.5，所以 p 也是一个随机变量，它服从Beta分布。如果我们要确定的不只是硬币这种 $\{0,1\}$ 二元取值的，而是骰子这种多元的，对应的事件成功的概率 p 服从狄利克雷分布

第四集：牛顿方法

- 广义线性模型

- 下面来看如何在指数分布族的基础上得到对应的线性模型，这些模型都属于广义线性模型。对于广义线性模型来说需要满足下面三个前提条件
 - $y | x; \theta \sim \text{ExponentialFamily}(\eta)$ ：说明 y 关于参数 η 的分布满足指数分布族的定义（其中 η 由自变量 x 和参数 θ 通过某种运算得到，该运算在广义线性模型中是有一定条件的，该条件详见下面第三条）
 - 广义线性模型的目标是在给定 x 的情况下需要得到 $T(y)$ 的期望，即 $E[T(y)|x]$ 的值（ $T(y)$ 为充分统计量，在统计学的角度看可以提供参数 θ 的全部信息，在通常情况下等于 y 值，因此我们相当于在给定 x 的情况下得到 $E[y|x]$ 的值，也就是原分布的期望值）
 - 对照我们原来的线性回归和Logistic回归，我们发现在线性回归中我们要预测的 $h(\theta)$ 就是对应正态分布的均值 $\theta^T x$ ，在Logistic回归中我们要预测的 $h(\theta)$ 就是 $E[y|x]$ 。由此线性回归和Logistic回归都是满足这个的定义的
 - $p(y = 1|x; \theta) = 0 \cdot p(y = 0|x; \theta) + 1 \cdot p(y = 1|x; \theta) = E[y|x; \theta]$
 - 广义线性模型的自然参数 η 与自变量 x 满足线性关系，即 $\eta = \theta^T x$
 - 如果 η 为向量，则 $\eta_i = \theta_i^T x$

第四集：牛顿方法

• 广义线性模型

- 满足上述三个条件的模型都可统称为广义线性模型，在广义线性模型中有很多好的性质（比如说学习起来速度快，比较高效）
- 下面来看之前的Logistic回归是怎样对应到广义线性模型上的
 - 先来看Logistic回归是如何满足广义线性模型的三个条件的

$$\begin{aligned} p(y; \phi) &= \phi^y (1 - \phi)^{1-y} \\ &= \exp(y \log \phi + (1 - y) \log(1 - \phi)) \\ &= \exp \left(\left(\log \left(\frac{\phi}{1 - \phi} \right) \right) y + \log(1 - \phi) \right) \end{aligned}$$

在之前的例子中我们已经看过Logistic回归对应的二项分布属于指数分布族；左下图展示了Logistic回归定义的 $h(x)$ 是符合 $E[y|x]$ 的；同时在Logistic回归中 $\eta = \theta^T x$ ，符合第三个条件

$T(y) = y$	$h_{\theta}(x) = E[y x; \theta]$
$a(\eta) = -\log(1 - \phi)$	$= \phi$
$= \log(1 + e^{\eta})$	$= 1/(1 + e^{-\eta})$
$b(y) = 1$	$= 1/(1 + e^{-\theta^T x})$

左下图用 η 代入其中的 ϕ 后得到的函数 $g(\eta)$ 称为正则响应函数(canonical response function)，而 g^{-1} 称为正则关联函数(canonical link function)

第四集：牛顿方法

- 广义线性模型

- 接下来看一下之前的线性回归是怎样对应到广义线性模型上的

$p(y; \mu) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(y - \mu)^2\right)$	$\eta = \mu$	$h_{\theta}(x) = E[y x; \theta]$
$= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}y^2\right) \cdot \exp\left(\mu y - \frac{1}{2}\mu^2\right)$	$T(y) = y$	$= \mu$
	$a(\eta) = \mu^2/2$	$= \eta$
	$= \eta^2/2$	$= \theta^T x.$
	$b(y) = (1/\sqrt{2\pi}) \exp(-y^2/2)$	

- 由此我们可以总结出机器学习的一个简单的自动化方法
 - 先根据数据的特征选择对应的符合指数分布族的分布，然后根据这个分布通过计算 $E[T(y)|x]$ 得到对应的式子 $f(\varphi)$ ，用 η 代入其中的 φ 后得到正则响应函数 $g(\eta)$ ，再将式子中的 η 替换为 $\theta^T x$ ，即可得到对应的广义线性模型 $h(x)$ 。接下来按照之前的办法写出对数似然方程，再通过梯度上升或者牛顿方法即可完成对于系数 θ 的求解

第四集：牛顿方法

- 广义线性模型

- 下面我们用这个思想来看一个复杂一点的例子，考虑一个多个离散变量的分类问题

- 首先看到这个是个分类问题，说明最后结果的取值不是连续的而是离散的；而条件又告诉我们，最后的取值不是二元的，因此不能用二项分布。由于是多元的变量，因此我们选定我们的模型符合多项式分布
 - 接下来我们需要将多项式分布写成指数分布族的形式
 - 首先先看多项式分布是什么样子的。我们可以假设有 k 个离散的结果为 $\{y_1, y_2, \dots, y_k\}$ ，而这些结果发生的概率分别为 $\varphi_1, \varphi_2, \dots, \varphi_{k-1}, 1 - (\varphi_1 + \varphi_2 + \dots + \varphi_{k-1})$ （避免过度参数化，多出一个多余参数 φ_k ），从而我们的参数集中只有 $k-1$ 个元素 $\{\varphi_1, \varphi_2, \dots, \varphi_{k-1}\}$ ，但是为了后面简单表示，仍然保留了 φ_k 这一写法。
 - 然后我们定义一个记号 $1\{\cdot\}$ ，在这个中括号里面是一个命题，当这个命题为真时该值为1，否则该值为0，即 $1\{\text{True}\} = 1, 1\{\text{False}\} = 0$ ，例子： $1\{2 = 3\} = 0, 1\{5 - 2 = 3\} = 1$

第四集：牛顿方法

- 广义线性模型

- 续上需要将多项式分布写成指数分布族的形式

- 借助上面的记号，我们可以将 $p(y; \phi)$ 表达成下面的通项形式

$$\begin{aligned} p(y; \phi) &= \phi_1^{1\{y=1\}} \phi_2^{1\{y=2\}} \dots \phi_k^{1\{y=k\}} \\ &= \phi_1^{1\{y=1\}} \phi_2^{1\{y=2\}} \dots \phi_k^{1 - \sum_{i=1}^{k-1} 1\{y=i\}} \end{aligned}$$

- 接下来的问题是怎样将 $1\{\cdot\}$ 表示成可解释的函数形式，因此我们可以引入一个 $(k-1) \times k$ 维的矩阵 M ，其定义如下（仅 m_{ii} 上的元素为1，否则为0（ $i=1 \sim k-1$ ））

$$M = [m(1) \quad m(2) \quad \dots \quad m(k-1) \quad m(k)] = \begin{bmatrix} 1 & 0 & \dots & 0 & 0 \\ 0 & 1 & & 0 & 0 \\ \vdots & & \ddots & \vdots & \\ 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & & 1 & 0 \end{bmatrix}$$

从而我们可以将该矩阵与这一记号结合起来，得到 $(m(y))_i = 1\{y = i\}$,

第四集：牛顿方法

- 广义线性模型

- 续上需要将多项式分布写成指数分布族的形式

$$\begin{aligned}P(y; \phi) &= \phi_1^{1\{y=1\}} \phi_2^{1\{y=2\}} \dots \phi_k^{1\{y=k\}} = \phi_1^{1\{y=1\}} \phi_2^{1\{y=2\}} \dots \phi_k^{1 - \sum_{i=1}^{k-1} 1\{y=i\}} = \phi_1^{(m(y))1} \phi_2^{(m(y))2} \dots \phi_k^{1 - \sum_{i=1}^{k-1} (m(y))i} \\&= \exp \left[(m(y))1 \log(\phi_1) + (m(y))2 \log(\phi_2) + \dots + \left(1 - \sum_{i=1}^{k-1} (m(y))i\right) \log(\phi_k) \right] \\&= \exp \left[(m(y))1 \log(\phi_1/\phi_k) + (m(y))2 \log(\phi_2/\phi_k) + \dots + (m(y))k-1 \log(\phi_{k-1}/\phi_k) + \log(\phi_k) \right] \\&= \exp \left[\begin{bmatrix} \log(\phi_1/\phi_k) & \log(\phi_2/\phi_k) & \dots & \log(\phi_{k-1}/\phi_k) \end{bmatrix} \begin{bmatrix} (m(y))1 \\ (m(y))2 \\ \vdots \\ (m(y))k-1 \end{bmatrix} + \log(\phi_k) \right]\end{aligned}$$

$$p(y; \eta) = b(y) \exp(\eta^T T(y) - a(\eta))$$

对比指数分布族可以发现 $m(y)$ 就是这里的 $T(y)$ ，并且可以得到如右图所示的关系

$$\begin{aligned}\eta &= \begin{bmatrix} \log(\phi_1/\phi_k) \\ \log(\phi_2/\phi_k) \\ \vdots \\ \log(\phi_{k-1}/\phi_k) \end{bmatrix} \\a(\eta) &= -\log(\phi_k) \\b(y) &= 1.\end{aligned}$$

第四集：牛顿方法

- 广义线性模型
 - 接下来我们需要计算 $E[T(y)|x]$

$$\begin{aligned}h_{\theta}(x) &= E[T(y)|x; \theta] \\&= E \left[\begin{array}{c} 1\{y = 1\} \\ 1\{y = 2\} \\ \vdots \\ 1\{y = k - 1\} \end{array} \middle| x; \theta \right] \\&= \begin{bmatrix} \phi_1 \\ \phi_2 \\ \vdots \\ \phi_{k-1} \end{bmatrix}\end{aligned}$$

这里面利用了前面定义的记号 $1\{\cdot\}$ 和性质 $(m(y))_i = 1\{y = i\}$ （由于在前面已经证实 $(m(y))_i = (T(y))_i$ ，因此这也就等价于 $(T(y))_i = 1\{y = i\}$ ），可以不难推出 $E[(T(y))_i] = P(y = i) = \phi_i$

接下来的问题就是将 ϕ 表示成 η 经过某种运算的结果，得到正则响应函数 $g(\eta)$

$$\eta = \begin{bmatrix} \log(\phi_1/\phi_k) \\ \log(\phi_2/\phi_k) \\ \vdots \\ \log(\phi_{k-1}/\phi_k) \end{bmatrix}$$

第四集：牛顿方法

- 广义线性模型

- 下面需要得到正则响应函数 $g(\eta)$

- 不失一般性地，我们可以作如下的假设

$$\eta_i = \log \frac{\phi_i}{\phi_k}$$

这里面的 $i=1,2,\dots,k$ ，此时我们可以看到 $\eta_k = \log 1 = 0$ ，不会影响结果的正确性

- 从而可以得到

$$\begin{aligned} e^{\eta_i} &= \frac{\phi_i}{\phi_k} \\ \phi_k e^{\eta_i} &= \phi_i \\ \phi_k \sum_{i=1}^k e^{\eta_i} &= \sum_{i=1}^k \phi_i = 1 \end{aligned} \quad \longrightarrow \quad \phi_i = \frac{e^{\eta_i}}{\sum_{j=1}^k e^{\eta_j}}$$

第四集：牛顿方法

- 广义线性模型

- 由此我们得到了正则响应函数，通过将式子中的 η 替换为 $\theta^T x$ ，即可得到对应的广义线性模型 $h(x)$ 和对数似然方程 $\ell(\theta)$

$$\begin{aligned} h_{\theta}(x) &= E[T(y)|x; \theta] \\ &= E \left[\begin{array}{c} 1\{y=1\} \\ 1\{y=2\} \\ \vdots \\ 1\{y=k-1\} \end{array} \middle| x; \theta \right] \\ &= \begin{bmatrix} \phi_1 \\ \phi_2 \\ \vdots \\ \phi_{k-1} \end{bmatrix} \\ &= \begin{bmatrix} \frac{\exp(\theta_1^T x)}{\sum_{j=1}^k \exp(\theta_j^T x)} \\ \frac{\exp(\theta_2^T x)}{\sum_{j=1}^k \exp(\theta_j^T x)} \\ \vdots \\ \frac{\exp(\theta_{k-1}^T x)}{\sum_{j=1}^k \exp(\theta_j^T x)} \end{bmatrix} \end{aligned}$$

在这里面只定义了 $\phi_1, \phi_2, \dots, \phi_{k-1}$ 的求解方法，而根据之前所说， ϕ_k 可以通过 $1 - (\phi_1 + \phi_2 + \dots + \phi_{k-1})$ 得到

$$\begin{aligned} \ell(\theta) &= \sum_{i=1}^m \log p(y^{(i)} | x^{(i)}; \theta) \\ &= \sum_{i=1}^m \log \prod_{l=1}^k \left(\frac{e^{\theta_l^T x^{(i)}}}{\sum_{j=1}^k e^{\theta_j^T x^{(i)}}} \right)^{1\{y^{(i)}=l\}} \end{aligned}$$

这种方法也被称为Softmax回归(Softmax Regression)，可以看成是Logistic回归的一种推广形式