

Linear and Non-linear Regression

created

2015-Jul-17

updated:

views: 198'

pages: 22

likes:

by Ramesh Hariharan

made with raindrops



The Regression Problem

- ullet Given an (unknown) function g(x) evaluated at n points $x_1 \ldots x_n$
- ullet The goal is recover $oldsymbol{g}$
- ullet Of course, $oldsymbol{g}$ is not specified uniquely by these $oldsymbol{n}$ evaluations
- ullet So assume g is of a simple, given form, say ax+b or ax^2+bx+c , where a,b,c are unknown parameters
- ullet Find the a,b,c that match the evaluations most closely

The Regression Problem More Precisely

- Given n pairs x_i, y_i
- $ullet x_i \in R^d$, $y_i \in R$
- ullet Given a model $f_{\Theta}(x)$, $x\in R^d$, Θ represents the unknown parameters
- ullet Find Θ that minimizes $h(\Theta) = \sum_i |f_\Theta(x_i) y_i|^2$

Linear Regression

- ullet Assume $f_{\Theta}(x) = x^T \Theta$, $\Theta, x \in R^d$
- ullet Find Θ that minimizes $h(\Theta) = \sum_i |x_i^T\Theta y_i|^2$
- Set $\frac{\partial h}{\partial \Theta}$ to 0
- Check that this is indeed a minimum, not a maximum, not a saddle etc

Determining Parameters in Linear Regression

- Rewriting more cleanly...
- $ullet h(\Theta) = (X\Theta y)^T (X\Theta y)$, where X is n imes d , the ith row of X is x_i
- ullet $=2X^T(X\Theta-y)$, where $rac{dh}{d\Theta}$ is a vector of partial derivatives w.r.t each of the entries in the vector Θ
- ullet Setting these to 0, we get a system of d equalities $X^TX\Theta=X^Ty$. Solution: $(X^TX)^{-1}X^Ty$

Global Minimality

- Invertibility? One solution if invertible, many otherwise. Why?
- Each solution of the form $\Theta = (X^T X)^{-1} X^T Y$ is a global minimum. Why?
- ullet By Taylor's expansion for multivariate functions, for any vector $\Delta=(\Delta_1,\Delta_2,\ldots,\Delta_d)\in R^d$,

$$h(\Theta + \Delta) = h(\Theta) + \Delta^T h'(\Theta) + rac{1}{2} \Delta^T h''(\Theta) \Delta$$

• Higher order terms are 0 for linear regression

Multivariate Taylor's Theorem Outline

- ullet r(t) = s(x(t),y(t)) Univariate Taylor: $r(t) = r(0) + tr'(0) + rac{t^2}{2}r''(0) + \cdots$
- $ullet r'(t) = rac{\partial s(x,y)}{\partial x} rac{\partial x}{\partial t} + rac{\partial s(x,y)}{\partial y} rac{\partial y}{\partial t}$
- $r''(t) = \frac{\partial^2 s(x,y)}{\partial x \partial x} (\frac{\partial x}{\partial t})^2 + 2 \frac{\partial^2 s(x,y)}{\partial y \partial x} \frac{\partial y}{\partial t} \frac{\partial x}{\partial t} + \frac{\partial^2 s(x,y)}{\partial y \partial y} (\frac{\partial y}{\partial t})^2$ plus terms which involve second derivative w.r.t t
- $ullet x(t) = x_0 + \Delta_1 t$, $y(t) = y_0 + \Delta_2 t$
- $ullet s(x_0 + \Delta_1 t, y_0 + \Delta_2 t) = s(x_0, y_0) + t \Delta^T s'(x_0, y_0) + rac{t^2}{2} \Delta^T s''(x_0, y_0) \Delta + \cdots$

Global Minimality

$$\bullet \ h'(\Theta) = \begin{pmatrix} \frac{\partial h}{\partial \Theta_1} \\ \frac{\partial h}{\partial \Theta_2} \\ \vdots \\ \frac{\partial h}{\partial \Theta_d} \end{pmatrix}, \ h''(\Theta) = \begin{pmatrix} \frac{\partial^2 h}{\partial \Theta_1 \partial \Theta_1} & \frac{\partial^2 h}{\partial \Theta_1 \partial \Theta_2} & \cdots & \frac{\partial^2 h}{\partial \Theta_1 \partial \Theta_d} \\ \frac{\partial^2 h}{\partial \Theta_2 \partial \Theta_1} & \frac{\partial^2 h}{\partial \Theta_2 \partial \Theta_2} & \cdots & \frac{\partial^2 h}{\partial \Theta_2 \partial \Theta_d} \\ \ddots & & & & \\ \frac{\partial^2 h}{\partial \Theta_d \partial \Theta_1} & \frac{\partial^2 h}{\partial \Theta_d \partial \Theta_2} & \cdots & \frac{\partial^2 h}{\partial \Theta_d \partial \Theta_d} \end{pmatrix}$$

- We chose Θ so $h'(\Theta)=0$.
- Claim: $\Delta^T h''(\Theta) \Delta \geq 0$, which suffices for global minimality

Global Minimality and Uniqueness

- $h''(\Theta)$, also called the Hessian, is exactly $2X^TX!$
- $\Delta^T X^T X \Delta = |X\Delta|^2 \ge 0$
- If X (which is $n \times d$, $n \geq d$) has full rank d then $X\Delta$ is always non-zero, so $|X\Delta|^2 > 0$, i.e., there is a unique global minimum
- Otherwise, any Δ in the nullspace of X satisfies $X\Delta=0$ and $h(\Theta+\Delta)=h(\Theta)$, so moving around in the nullspace yields all global minima

Invertibility of X^TX

- $rank(X) = rank(X^TX)$. Why? Because the nullspaces are the same.
- ullet If $oldsymbol{X}$ has rank $oldsymbol{d}$ then $oldsymbol{X^TX}$ is invertible and there is an unique global minimum
- ullet If X has smaller rank, then X^TX is not invertible, then how do we find solutions to $X^TX\Theta=X^Ty$?
- ullet Note X^TX takes an arbitrary vector $oldsymbol{\Theta}$ and transforms it into the column space of X^T .

The Pseudo-Inverse of $oldsymbol{X^TX}$

- ullet Since X^TX is a real, symmetric, positive semi-definite d imes d matrix, all its eigenvalues are real and non-negative
- ullet The corresponding eigenvectors are either in the column space of X^T or orthogonal to it, and form an orthonormal basis
- ullet Eigenvector $oldsymbol{v}$ has eigenvalue 0 if it is orthogonal to all the columns of $oldsymbol{X^T}$ and eigenvalue non-0 if it is in the column space of $oldsymbol{X^T}$
- ullet $X^TX = \sum_i e_i v_i v_i^T$ where e_i are the non-0 eigenvalues and v_i denotes the corresponding eigenvectors

The Pseudo-Inverse of $oldsymbol{X^TX}$ Contd.

- ullet Consider $(\sum_i rac{1}{e_i} v_i v_i^T) X^T y$ as a solution for Θ
- ullet Claim: $(X^TX)(\sum_i rac{1}{e_i} v_i v_i^T) X^T y = X^T y$
- Proof:
 - $egin{array}{l} egin{array}{l} igl(\sum_i e_i v_i v_i^T) (\sum_i rac{1}{e_i} v_i v_i^T) X^T y = (\sum_i v_i v_i^T) X^T y \end{array}$
 - $(\sum_i v_i v_i^T) x = x$ for any vector x in the column space of X^T
 - \circ Since X^Ty is in the column space of X^T , $(\sum_i v_i v_i^T) X^Ty = X^Ty$, as needed

Summarizing the Non-Invertible Case

- Since the column space of X^T has dimension at least 1, at least one eigenvalue of X^TX is non-zero, hence there is at least one solution satisfying $h'(\Theta)=0$, namely $\Theta=(\sum_i \frac{1}{e_i} v_i v_i^T) X^T y$
- ullet This is a global minimum, though not unique, any $h(\Theta+\Delta)$ where Δ is in the nullspace of X (equivalently, orthogonal to column space of X^T) is also a global minimum
- ullet There is a unique global minimum in the column space of X^T

Running Time

- ullet Compute X^TX
- Compute eigenvalues and eigenvectors
- ullet Compute $(\sum_i rac{1}{e_i} v_i v_i^T) X^T y$
- Naively, $O(d^2n + d^3)$
- ullet There are ways to compute $oldsymbol{X^TX}$ faster, approximately, via random sampling

The Non-Linear Case

- Say $f_{\Theta}(x)=(x^T\Theta)^2$ but not $f_{\Theta}(x)=(x^2)^T\Theta$, where x^2 denotes the vector x with entries squared; assume well-behaved (continuously differentiable)
- $h(\Theta + \Delta) = h(\Theta) + \Delta^T h'(\Theta) + \Delta^T h''(\Theta) \Delta + O(\Delta^3)$ for small enough Δ
- $h'(\Theta)=0$ can no longer be solved using linear equations. Needs a general root-finding procedure
- $h''(\Theta)$ need not be of the form A^TA , which means $\Delta^Th''(\Theta)\Delta$ could be negative, i.e., a root need not be local minimum, let alone a global one

The Gauss-Newton Approach

- ullet Start with a particular Θ_0
- ullet Linearize $f_{\Theta}(x)$ in the neighborhood of this Θ_0 : $f_{\Theta_0+\Delta}(x)\sim f_{\Theta_0}(x)+[f'_{\Theta_0}(x)]^T\Delta$
- ullet Minimize $hh(\Delta) = \sum_i |f_{\Theta_0}(x_i) + [f'_{\Theta_0}(x_i)]^T \Delta y_i|^2$
- ullet Use linear regression to find $\Delta=(Z^TZ)^{-1}Z^T(y-f_{\Theta_0}(x))$ that minimizes $hh(\Delta)$; iterate with $\Theta_0=\Theta_0+\Delta$
- ullet Where the ith row of Z is the vector $[f_{\Theta_0}'(x_i)]^T$

Convergence of the Gauss-Newton Approach

- If $f_{\Theta_0+\Delta}(x)\sim f_{\Theta_0}(x)+[f'_{\Theta_0}(x)]^T\Delta$ is *close* to an equality then $\Theta_0+\Delta$ provides a strictly better fit than Θ_0 unless Θ_0 is already a local minimum
- So if we start very close to a local minimum (so the linear approximation holds in a very close neighborhood), this will converge
- But no guarantee if you start further away.

Modified Gauss-Newton with Smaller Step Sizes

- $\Delta=(Z^TZ)^{-1}Z^T(y-f_{\Theta_0}(x))$, which equals $-rac{1}{2}(Z^TZ)^{-1}h'(\Theta_0)$ (why?)
- $ullet h(\Theta) = (f_{\Theta}(x) y)^T (f_{\Theta}(x) y)$ and $h'(\Theta) = 2Z^T (f_{\Theta}(x) y)$
- ullet $\Theta_0 + lpha \Delta = \Theta_0 rac{lpha}{2} (Z^T Z)^{-1} h'(\Theta_0)$
- For small enough $\alpha>0$, $h(\Theta_0+\alpha\Delta)\sim h(\Theta_0)-\frac{\alpha}{2}h'(\Theta_0)^T(Z^TZ)^{-1}h'(\Theta_0)< h(\Theta_0)$, at least if Z is full rank and Θ_0 is not already a local minimum

The Modified Gauss-Newton Algorithm for the Non-Linear Case

- ullet Start with some Θ_0 and then move to $\Theta_0-lpha(Z^TZ)^{-1}h'(\Theta_0)=\Theta_0-lpha(Z^TZ)^{-1}Z^T(f_{\Theta_0}(x)-y)$
- ullet Here the ith row of Z is the vector $[f_{\Theta_0}{}'(x_i)]^T$
- ullet lpha can be determined by single-parameter minimization
 - e.g., bisection: start with an interval, evaluate derivative in the middle, if negative, go right,
 otherwise go left
- Stop when there isn't much progress

Some More Modifications: Levenberg

- ullet Instead of $\Theta_0-lpha(Z^TZ)^{-1}h'(\Theta_0)$ one could simply follow steepest gradient to $\Theta_0-lpha Ih'(\Theta_0)$
- ullet The latter always leads to improvement but could be slow, the former is quicker if $f_{\Theta_0}(x)$ is linear in the neighbourhood of Θ_0
- ullet Levenberg: Combine the two approaches, i.e., use $\Theta_0 lpha(Z^TZ + \lambda I)^{-1}h'(\Theta_0)$
- ullet Adjust $oldsymbol{\lambda}$ in each iteration, if there is quick descent then keep it low, otherwise increase it so you are closer to steepest descent

Some More Modifications: Marquardt

ullet Marquardt: Instead of $\Theta_0-lpha(Z^TZ+\lambda I)^{-1}h'(\Theta_0)$, use $\Theta_0-lpha(Z^TZ+\lambda\mathrm{diag}(Z^TZ))^{-1}h'(\Theta_0)$

$$\bullet \ \mathrm{diag}(Z^TZ) = \begin{pmatrix} \sum_i (\frac{\partial f_{\Theta_0}(x_i)}{\partial \Theta_1})^2 & & & \\ & \sum_i (\frac{\partial f_{\Theta_0}(x_i)}{\partial \Theta_2})^2 & & & \\ & & \ddots & & \\ & & & \sum_i (\frac{\partial f_{\Theta_0}(x_i)}{\partial \Theta_d})^2 \end{pmatrix}$$

• When λ is large so one is in the steepest descent regime, move faster in directions where the partial derivative is smaller, in a bid to accelerate

Finishing Up

• Can we use these regression methods, particularly Levenberg-Marquardt, to identify the orbit of Mars, based on data in the previous lecture?