

# Lecture 8 - scribbles

Friday, September 28, 2018 13:34

today: - logistic regression  
- numerical optimization  
IRLS

## Logistic regression model:

$$Y = \{0, 1\}$$
$$p(y=1|x) = \sigma(w^T x)$$

$$p(y=0|x) = 1 - \sigma(w^T x) = \sigma(-w^T x)$$

ie.  $Y|X=x$  is Bernoulli( $\sigma(w^T x)$ )

$$p(y|x) = \sigma(w^T x)^y \sigma(-w^T x)^{1-y}$$

given  $(x_i, y_i)_{i=1}^n$ , maximum conditional likelihood:

$$l(w) = \sum_{i=1}^n \log p(y_i|x_i; w) = \sum_{i=1}^n [y_i \log \sigma(w^T x_i) + (1-y_i) \log \sigma(-w^T x_i)]$$

derivative

$$\sigma'(z) = \sigma(z)\sigma(-z)$$

$$\nabla_w \sigma(w^T x_i) = x_i [\sigma(w^T x_i)\sigma(-w^T x_i)] \quad \text{let } v_i \triangleq w^T x_i$$

$$\nabla l(w) = \sum_{i=1}^n x_i \left[ y_i \frac{\sigma(v_i)\sigma(-v_i)}{\sigma(v_i)} - (1-y_i) \frac{\sigma(v_i)\sigma(-v_i)}{\sigma(-v_i)} \right]$$

$$\left[ y_i [\underbrace{\sigma(-v_i) + \sigma(v_i)}] - \sigma(v_i) \right]$$

$$[ \text{if } Y = \{\pm 1\}$$

$$\text{encode } p(y|x) = \sigma(yw^T x)$$

$$\nabla \ell(w) = \sum_{i=1}^n x_i [y_i - \sigma(w^T x_i)]$$

solve for  $\nabla \ell(w) = 0 \Rightarrow$  need to solve a transcendental eq.

$\Downarrow$  need to use numerical methods

because  $\frac{1}{1 + \exp(-w^T x_i)} (\dots) = 0$   
 $\uparrow$  need to solve for this

contrast to least square regression.

$$\nabla \ell(w) = \sum_{i=1}^n x_i [y_i - w^T x_i]$$

$\uparrow$   
linear in  $w$

numerical optimization

want to minimize  $f(w)$

$$\begin{array}{ll} \min & f(w) \\ \text{s.t.} & w \in \mathbb{R}^d \end{array}$$

1) gradient descent (1st order method)

start at  $w_0$

$$w_{t+1} = w_t - \underbrace{\alpha_t}_{\text{step-size}} \nabla f(w_t)$$

step-size rules:

a) constant step-size  $\gamma_t = \frac{1}{L}$  ← Lipschitz constant of  $\nabla f$

$$\|\nabla f(w) - \nabla f(w')\| \leq L \|w - w'\|$$

b) decreasing step-size rule:  $\gamma_t = \frac{c}{t}$  ← constant

usually want:  $\sum_t \gamma_t = +\infty$      $\sum_t \gamma_t^2 < \infty$

c) choose  $\gamma_t$  by "line search" :  $\min_{\gamma \in \mathbb{R}} f(w_t + \gamma d_t)$

direction for update  
(e.g.  $-\nabla f(w_t)$ )

↙ costly in general

instead do approximate search  
e.g. Armijo line search  
(see Boyd's book)

2) Newton's method (2nd order method)

motivation: minimizing a quadratic approximation:

taylor expansion:  $f(w) = f(w_t) + \nabla f(w_t)^T (w - w_t) + \frac{1}{2} (w - w_t)^T H(w_t) (w - w_t)$

Hessian  $[H(w_t)]_{ij} = \frac{\partial^2 f(w_t)}{\partial w_i \partial w_j}$

+  $O(\|w - w_t\|^3)$

↑  
Taylor's remainder

$$= \underbrace{Q_t(w)} + O(\|w - w_t\|^3)$$

quadratic model approximation

$w_{t+1} \rightarrow$  minimizing  $Q_t(w)$

$$\nabla_w Q_t(w) = 0$$

$$\nabla f(w_t) + H(w_t)(w - w_t) = 0$$

$$\Rightarrow w - w_t = -H^{-1}(w_t) \nabla f(w_t)$$

$$w_{t+1} = w_t - \underbrace{H^{-1}(w_t)}_{\text{inverse Hessian}} \nabla f(w_t)$$

Newton's update

inverse Hessian  $\rightarrow O(d^3)$  time to compute in general  $O(d^2)$  space

Damped Newton: you add a stepsize to stabilize Newton's method

$$w_{t+1} = w_t - \underbrace{\gamma}_{\text{step-size}} H^{-1}(w_t) \nabla f(w_t)$$

why Newton's method?

- much faster convergence in # of iterations vs gradient descent
- affine invariant  $\rightarrow$  invariant to rescaling of variables

but iterations are costly  $O(d^3)$  time  $O(d^2)$  memory vs.  $O(d)$  for GD.

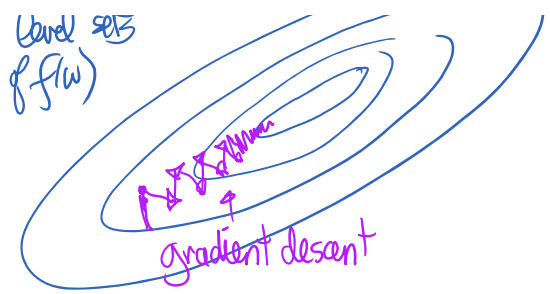
level sets of  $f(w)$



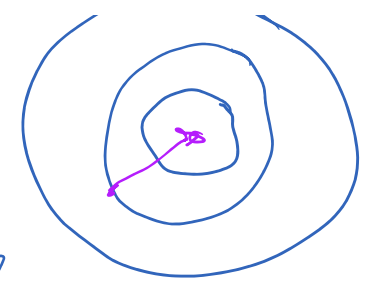
Newton is transforming space

using Hessian to make it "round"





using Hessian to make it "well-conditioned"



$$z_i = H^{-1}x$$

level sets of  $f(z)$

$$\nabla f(z) = H^{-1} \nabla f(x)$$

$$\frac{1}{2} x^T H x = c$$

$$\frac{1}{2} x^T P^T \Sigma P x = c$$

look at "quadratic forms"

(side note on implicit regularization properties of optimization method)

Newton's method for logistic regression: IRLS

recall for log. reg.,  $\nabla l(w) = \sum_{i=1}^n x_i [y_i - \sigma(w^T x_i)]$

$$H(l(w)) = \sum_{i=1}^n x_i x_i^T \sigma(w^T x_i) \sigma(-w^T x_i)$$

$$v^T H v = - \sum_{i=1}^n \underbrace{(v^T x_i)(x_i^T v)}_{(v^T x_i)^2 \geq 0} \underbrace{\sigma(-) \sigma(-)}_{\geq 0}$$

ie.  $v^T H v \leq 0 \forall v \in \mathbb{R}^d$   
 $\Rightarrow H \preceq 0$   
 ie. concave function

notation:

recall  $X = \begin{pmatrix} - & \vdots & - \\ - & x_i^T & - \\ & \vdots & \\ - & \vdots & - \end{pmatrix}$

let  $\mu_i \triangleq \sigma(w^T x_i) \in ]0, 1[$

Let  $\mu_i \triangleq \sigma(w^T x_i) \in ]0,1[$

$$\nabla \ell(w) = \sum_{i=1}^n x_i [y_i - \mu_i] = X^T (y - \mu)$$

$$\text{Hessian: } = - \sum_{i=1}^n x_i x_i^T \mu_i (1 - \mu_i) = - X^T D X \quad \text{where } D_{ii} = \mu_i (1 - \mu_i)$$

[note: D depends]

Newton's update (here Newton's updates are maximizing log-likeli. because it is concave)  $\rightarrow$

Newton's update:

$$w_{k+1} = w_k - (-X^T D_k X)^{-1} X^T (y - \mu_k)$$

$$= (X^T D_k X)^{-1} [ (X^T D_k X) w_k + X^T (y - \mu_k) ]$$

$$\boxed{w_{k+1} = (X^T D_k X)^{-1} [ X^T D_k z_k ]} \quad \text{where } \boxed{z_k \triangleq X w_k + D_k^{-1} (y - \mu_k)}$$

this is a solution to "weighted least square problem"

$$\min_w \| D^{1/2} (z_k - Xw) \|^2$$

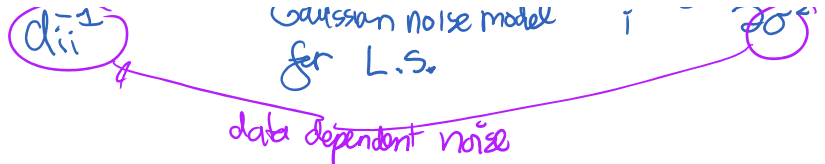
new target

$$\sum_{i=1}^n \frac{(z_i - x_i^T w)^2}{d_{ii}}$$

compare with Gaussian noise model for L.S.

$$\sum_i \frac{(y_i - x_i^T w)^2}{\sigma_i^2}$$

1-1



Newton's method for logistic regression  
 = iterative reweighted least square (IRLS)

2 comments for assignment:

- 1) stopping criterion:  $\|\nabla f(w_t)\| \leq \epsilon$
- 2) to compute  $A^{-1}V$       $A \setminus V$

Big data logistic regression:

suppose  $f(w) = \frac{1}{n} \sum_{i=1}^n f_i(w)$

- cannot do  $O(d^2)$  or  $O(d^3)$  operations  $\Rightarrow$  first order methods
- if  $n$  is huge, you cannot do batch method

$\nabla f(w) = \left( \sum_{i=1}^n \nabla f_i(w) \right)$   $O(nd)$  time

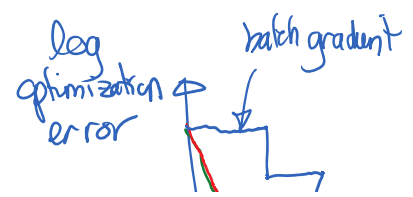
gradient of one  $f_i$  (pointing to  $\nabla f_i(w)$ )  
 "batch" gradient (pointing to the sum)

instead you "incremental gradient methods"

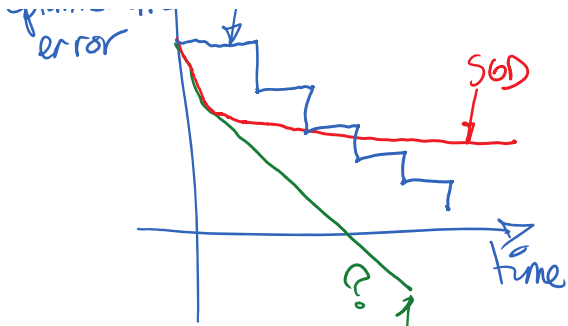
e.g. stochastic gradient descent (SGD):  $w_{t+1} = w_t - \gamma_t \nabla f_{i_t}(w_t)$   $O(d)$  time  
 where  $i_t$  is picked randomly

SGD  $\rightarrow$  cheap updates, but slower convergence per iteration

batch gradient  $\rightarrow$  expensive, but fast " " "



SGD



yes? variance reduced methods

→ SAG: stochastic averaged gradient  
[2012]

$$\text{G.D.} : w_{t+1} = w_t - \delta_t \frac{1}{n} \sum_{i=1}^n \nabla f_i(w_t)$$

$$\text{SAG} : w_{t+1} = w_t - \delta_t \frac{1}{n} \sum_{i=1}^n v_i \quad \text{memory}$$

$$\text{where } v_i = \nabla f_i(w_{\text{old}(i)})$$

at each  $t$ , update only one  $v_{i_t} \triangleq \nabla f_{i_t}(w_t)$

$$\text{SAGA} : w_{t+1} = w_t - \delta \left( \nabla f_{i_t}(w_t) + \underbrace{\frac{1}{n} \sum_{i=1}^n v_i - v_{i_t}}_{\text{variance reducing correction}} \right)$$

(default method for logistic regression in Scikit-learn)