

## Lecture 8 - logistic regression & optimization

Monday, September 30, 2024 12:33 PM

today : - logistic regression  
• numerical optimization

logistic regression:

setup : binary classification  $\gamma \in \{0, 1\}$   $x \in \mathbb{R}^d$

generative model motivation:

suppose only assumption is there exists pdf (densities) in  $\mathbb{R}^d$   
for each class conditional

$$\begin{aligned} p(x|Y=1) &\quad \{ \quad p(x|Y=0) \\ p(Y=1|x=x) &\approx \frac{p(Y=1, x=x)}{p(Y=1, x=x) + p(Y=0, x=x)} \quad \{ \quad p(x=x) \\ &= \frac{1}{1 + \frac{p(Y=0, x=x)}{p(Y=1, x=x)}} \approx \frac{1}{1 + \exp(-f(x))} \end{aligned}$$

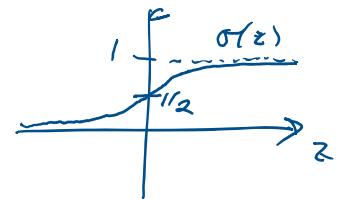
$$\text{where } f(x) \triangleq \log \frac{p(x=x | Y=1)}{p(x=x | Y=0)} \leftarrow \log \frac{p(Y=1)}{p(Y=0)}$$

"log odd" ratio      class conditional ratio      prior odd ratio

in general

$$p(Y=1|X=x) = \sigma(f(x)) \quad \text{where } \sigma(z) \triangleq \frac{1}{1 + \exp(-z)}$$

"sigmoid function"



$$\text{some properties of } \sigma(z): \quad \sigma(-z) = 1 - \sigma(z) \quad [\sigma(z) + \sigma(-z) = 1]$$

$$\frac{d\sigma(z)}{dz} = \sigma(z)\sigma(-z) = \sigma(z)(1-\sigma(z))$$

② to motivate linear logistic regression

consider class conditional in the exponential family

$$p(x|m) \triangleq h(x) \exp(M^T T(x) - A(m))$$

"canonical parameters"

Scalar fact.  $\rightarrow$  bay partition fact.  
 $\rightarrow$  normalize

"sufficient statistics"  
these specify the "flat" exponential family

Gaussian:

$$\log p(x|\mu, \sigma^2) = -\frac{1}{2} \log 2\pi\sigma^2 - \frac{(x-\mu)^2}{2\sigma^2}$$

$$= -\left( \frac{x^2}{2\sigma^2} - \frac{2x\mu}{\sigma^2} + \frac{\mu^2}{2\sigma^2} \right)$$

$\downarrow$   
 $\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$

$$\text{let } T(x) = \begin{bmatrix} -x^2/2 \\ x \end{bmatrix}$$

$$\mathcal{N}(\mu, \sigma^2) = \begin{bmatrix} 1/\sigma^2 \\ \mu/\sigma^2 \end{bmatrix}$$

$$A(\eta) = \frac{1}{2} \log(2\pi\sigma^2) + \frac{\mu^2}{2\sigma^2}$$

$$p(x|Y=1) = p(x|\eta_1)$$

$$p(x|Y=0) = p(x|\eta_0)$$

$$\text{log odds } f(x) = \log \frac{p(x|\eta_1)}{p(x|\eta_0)} + \log \frac{\frac{\pi}{1-\pi}}{1-\pi}$$

$$= (\eta_1 - \eta_0)^T T(x) + A(\eta_1) - A(\eta_0) + \log\left(\frac{\pi}{1-\pi}\right)$$

$$\triangleq w^T \varphi(x)$$

$$\text{where } w = \begin{pmatrix} \eta_1 - \eta_0 \\ A(\eta_1) - A(\eta_0) + \log\left(\frac{\pi}{1-\pi}\right) \end{pmatrix} \quad \varphi(x) = \begin{pmatrix} T(x) \\ 1 \end{pmatrix}$$

get logistic regression model

$$f_w(Y=1|x) = \sigma(w^T \varphi(x))$$

↓  
"feature map"

decision boundary

$$\{w^T \varphi(x) > 0\}$$

exercise for the reader : try argument with  $p(x|y) = N(x|\mu_y, \Sigma_y)$   
(hwk 2)

• if  $\Sigma_y = \Sigma_1$ , then  $\varphi(x) = \begin{pmatrix} x \\ 1 \end{pmatrix} \rightarrow$  "linear boundary in  $x^y$ "

• otherwise,  $\varphi(x) = \begin{pmatrix} -xx^T \\ x \\ 1 \end{pmatrix} \rightarrow$  quadratic decision boundary

## Linear Logistic regression model

$$p(y=1|x, w) = \sigma(w^T x) \quad y = \xi \{ 0, 1 \} \quad \text{decision rule: } \{ \xi w^T x > 0 \}$$

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

[if  $y = \xi \{ 1 \}$  "Rademacher R.V."]

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

$y|X=x$  is Bernoulli( $\sigma(w^T x)$ )

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

given  $(x_i, y_i)_{i=1}^n$  iid

Max conditional log likelihood to estimate  $w_{\text{acc}}$

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

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

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

$$\boxed{\nabla l(w) = \sum_{i=1}^n x_i \left[ y_i \left[ \underbrace{\sigma(-v_i) + \sigma(v_i)}_1 \right] - \sigma(v_i) \right]}$$

solve for  $\nabla l(w) = 0 \rightarrow$  need to solve transcendental equation

because  $\sum_i \frac{1}{1 + \exp(-w^T x_i)} (-) = 0$   
 ↗ need to use numerical methods  
 ↗ solve for  $w$

contrast to linear regression  $\nabla l(w) = \sum_{i=1}^n x_i [y_i - w^T x_i]$   
 ↗ linear in  $w$

13h28

## Numerical optimization

want to minimize  $f(w)$  (unconstrained) {suppose  $f$  is differentiable}

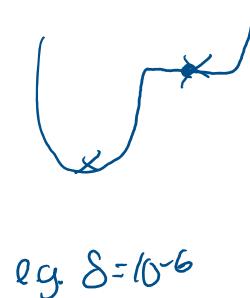
## Convex optimization

want to minimize  $f(w)$  (unconstrained) [suppose  $f$  is differentiable]  
 s.t.  $w \in \mathbb{R}^d$

### 1) gradient descent (1<sup>st</sup> order method)

start  $w_0$   
 iterate:  $w_{t+1} = w_t - \gamma_t \nabla f(w_t)$

stopping criterion: if  $\|\nabla f(w_t)\|^2 < \delta$   
 then stop



note: if  $f$  is  $\mu$ -strongly convex  $\Rightarrow f(w_t) - \min_w f(w) \leq \frac{1}{2\mu} \|\nabla f(w_t)\|^2$   
 $(\Leftrightarrow f(w) - \frac{\mu}{2} \|w\|^2 \text{ is convex})$

$$\text{if } \|\nabla f(w_t)\|^2 \leq \delta \Rightarrow \text{subopt.}(w_t) \leq \frac{\delta}{2\mu}$$

#### Step-size rules

a) constant step-size:  $\gamma_t = \frac{1}{L}$  [Lipschitz continuity constant for  $\nabla f$ ]

$$\text{i.e. } \|\nabla f(w) - \nabla f(w')\|_2 \leq L \|w - w'\|_2 \quad \forall w, w' \in \mathbb{R}^d$$

b) decreasing step-size rule

[this is more common for stochastic optimization]

$$\gamma_t = \frac{C}{t} \text{ constant}$$

$$\text{usually want } \sum_t \gamma_t = \infty$$

$$\sum_t \gamma_t^2 < \infty$$

$$\text{e.g. } f(w) \triangleq \mathbb{E}_{z \sim p} g(w, z)$$

$$w_{t+1} = w_t - \gamma_t \nabla_w g(w, z_t)$$

$$\text{e.g. } g(w, z) = f(y_i, h_w(z_i)) \text{ for ERM}$$

$$z = (x_i, y_i)$$

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

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

costly in general

→ instead do approximate search

e.g. Armijo line search [see Boyd's book]

#### Newton's method (2<sup>nd</sup> order method)

motivation: minimizing a quadratic approximation

↙ Hessian  
 $[\nabla^2 f(w)]_{ij} = \frac{\partial^2 f(w)}{\partial w_i \partial w_j}$

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

method is minimizing a quadratic approximation

Taylor expansion  
of  $w_t$

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

$\triangleq Q_t(w)$

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

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

Taylor's remainder

$w_{t+1} \rightsquigarrow$  obtain by minimizing  $Q_t(w)$

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

want

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

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

$$w_{t+1} = w_t - H^{-1}(w_t) \nabla f(w_t)$$

Newton's update

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

invert Hessian

$\rightsquigarrow O(d^3)$  time  
and  $O(d^2)$  space

$$H d_t = \nabla f$$

Note for hwk 2: solve for  $d_t$  in  $H d_t = \nabla f(w_t)$

$$\min_d \|H d - \nabla f(w_t)\|^2$$

use [numpy.linalg.lstsq] to solve this

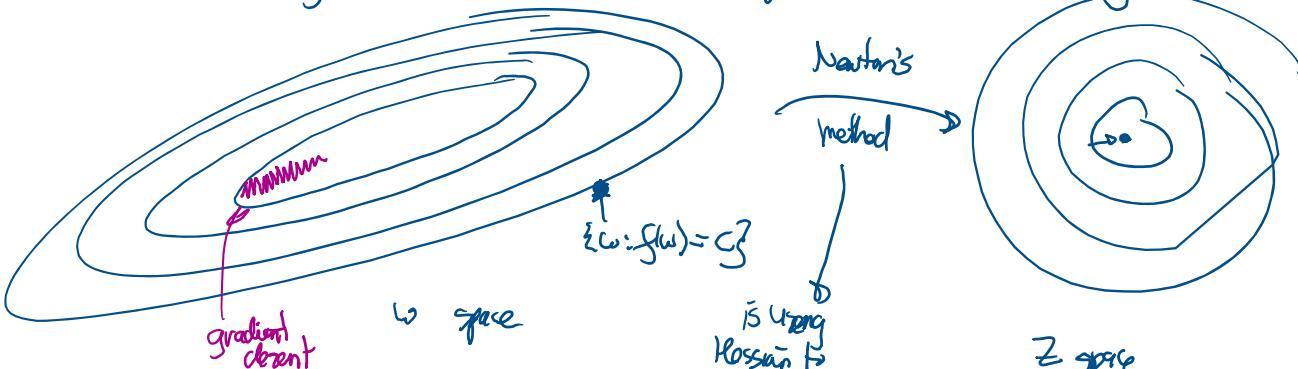
$\rightsquigarrow$  much more numerically stable

damped Newton: You add a step size to stabilize Newton's method

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

why Newton's method?

- much faster convergence in # iterations vs. gradient descent
- affine invariant  $\Rightarrow$  method's convergence is invariant to rescaling of variables



gradient descent

$\omega$  space

$z$  space

$\frac{1}{2} \omega^T H \omega = C$  augmented

 $\downarrow$ 
 $H = P^T Z P$ 

$(H$  is symmetric psd)

$\frac{1}{2} \omega^T P^T \xi P \omega = C \Leftrightarrow \frac{1}{2} z^T z = C$

$z \equiv Z^{1/2} \omega$

Exercise to reader :  $z_{t+1} = z_t - \gamma \nabla f(z_t)$

$\Leftarrow$

$w_{t+1} = w_t - \gamma H^{-1} \nabla f(w_t)$

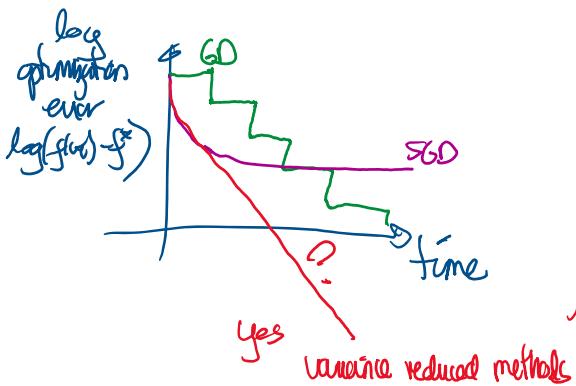
### Big data Logistic regression

- big  $d \Rightarrow$  cannot do  $O(d^2)$  or  $O(d^3)$  operations  $\Rightarrow$  first order methods
- if  $n$  is large, you cannot use batch methods  $\nabla f(\omega) = \left( \sum_{i=1}^n \nabla f_i(\omega) \right)$  gradient of loss on some data point  $i$

instead you can use "incremental gradient method"

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

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



### SAG & stochastic average gradient

$$GD: w_{t+1} = w_t - \gamma \sum_{i=1}^n \nabla f_i(w_t)$$

$$SAG [2012] \quad w_{t+1} = w_t - \gamma_t \sum_{i=1}^n v_i \quad \text{Via memory}$$

where  $v_i = \nabla f_i(w_{t+1})$

at each  $t$ , update  $v_{i_t} = \nabla f_i(w_t)$

$$w_{t+1} = w_t - \gamma_t (\nabla f_{i_t}(w_t) + \frac{1}{n} \sum_{j=1}^n (v_j - v_{i_t}))$$

(default method  
for large logistic regression)  
in Scikit-Learn

$\frac{1}{n}$  variance reduction come

# SVRG

Skipped part: Newton's method on logistic regression = IRLS:

Newton's method for logistic regression: IRLS

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

$$H(l(w)) = -\sum_{i=1}^n x_i x_i^T \sigma(w^T x_i) \sigma'(w^T x_i)$$

$$\nabla^T H \nabla = -\sum_{i=1}^n (\nabla x_i) (x_i^T \nabla) \underbrace{\sigma(-)}_{(x_i^T \nabla)^2 \geq 0} \underbrace{\sigma'(-)}_{\geq 0}$$

$\nabla^T H \nabla \leq 0$  Hverg

i.e. HJO  
i.e. concave fn.

Notation: recall

$$X = \begin{pmatrix} -x_1^T & \dots \\ \vdots & \end{pmatrix}_{n \times d}$$

Newton is maximizing  
instead of minimizing

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

$$\nabla l(w) = \sum_i x_i [y_i - \mu_i] = X^T (y - \mu)$$

$$\text{Hessian} = -\sum_i x_i x_i^T \mu_i (1 - \mu_i) = -X^T D(w) X$$

$$D \triangleq \begin{pmatrix} \vdots & \\ \sigma(w^T x_i) & \\ \vdots & \end{pmatrix}$$

$D_{ii} \triangleq \mu_i (1 - \mu_i)$   
diagonal matrix

$$\text{Newton's update: } w_{t+1} = w_t - (X^T D_t X)^{-1} X^T (y - \mu_t)$$

$$= (X^T D_t X)^{-1} \left[ (X^T D_t X) w_t + X^T (y - \mu_t) \right]$$

$$w_{t+1} = (X^T D_t X)^{-1} [X^T D_t z_t]$$

$$\text{where } z_t \triangleq X w_t + D_t^{-1} (y - \mu_t)$$

This is a solution to a "weighted least square problem"

$$\min_w \|D^{1/2} (z_t - X w)\|_2^2$$

$P$  new target  
 $w$  weights

$$\rightarrow \sum_i \frac{(z_i - w^T x_i)^2}{d_{ii}}$$

compare with Gaussian noise  
model for least square

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

Newton's method for logistic regression  
= iterated reweighted least squares (IRLS)