

today: • logistic regression  
• numerical optimization

logistic regression model:

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

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

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

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

given  $(x_i, y_i)_{i=1}^n$ , max. conditional log-likelihood to estimate  $\hat{w}_{ML}$

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

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

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

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

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

need to use numerical methods

solve for  $\nabla l(w) = 0 \Rightarrow$  need to solve a transcendental equation because  $\frac{1}{1+\exp(-v_i)}(\dots) = 0$

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

Numerical optimization:

want to minimize  $f(w)$  (unconstrained)  
st.  $w \in \mathbb{R}^d$

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

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



$w_0$   
 Iterate:  $w_{t+1} = w_t - \gamma_t \nabla f(w_t)$   
 Stopping criterion:  $\|\nabla f(w_t)\|^2 < \delta$  e.g.  $\delta = 10^{-6}$   
 Note:  $f$   $\mu$ -strongly convex  $\Rightarrow \|\nabla f(w_t)\|^2 < \delta \Rightarrow f(w_t) - f(w^*) \leq \frac{\delta}{2\mu}$   
 $\Leftrightarrow \mu > 0$   
 $f(w) - \frac{\mu}{2} \|w\|^2$  is convex

step-size rules:

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

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

b) decreasing step-size rule [this is more common for stochastic optimization]

e.g.  $f(w) \triangleq \mathbb{E}_{\xi} g(w, \xi)$

$\gamma_t = \frac{c}{t}$  constant

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

$\xi = \{x_i, y_i\}$

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

e.g.  $g(w, \xi) = \mathcal{L}(y_i, h(w; x_i))$   
 $\xi = (x_i, y_i)$

for ERM

c) choose  $\gamma_t$  by "line search":

$\min_{\gamma \in \mathbb{R}} f(w_t + \gamma \frac{dw_t}{dt})$

costly in general

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

\* instead do approximate search

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

15h26

Newton's method (2nd order method)

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

motivation: minimizing a quadratic approximation

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

+  $O(\|w - w_t\|^3)$   
 Taylor's remainder

$= \mathcal{Q}_t(w) + O(\|w - w_t\|^3)$

↑ Quadratic model app.

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

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

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

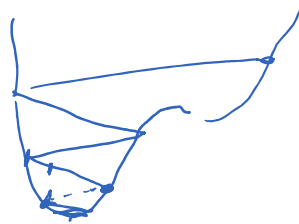
$$\Rightarrow w - w_t = -H^{-1}(w_t) \nabla f(w_t)$$

$$w_{t+1} = w_t - H^{-1}(w_t) \nabla f(w_t) \quad \text{Newton's update}$$

↳ Inverse Hessian  $\Leftrightarrow H^{-1} \nabla f$   
 $\rightarrow O(d^3)$  time to compute in general and  $O(d^2)$  space

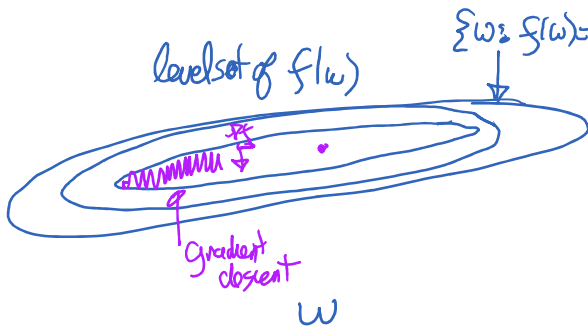
Damped Newton: you add a step-size to stabilize Newton's method

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

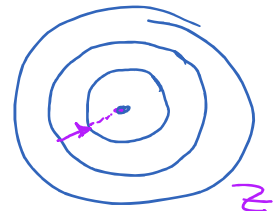


Why Newton's method?

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



Newton's method  
 is using Hessian to make f "well conditioned"  
 $z = H^{1/2} w$



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

↑ "quadratic form"

$$\frac{1}{2} w^T P^T \Sigma P w = c \quad \Leftrightarrow \frac{1}{2} z^T z = c$$

↑ diagonal

$$H = P^T \Sigma P$$

H is symmetric & PSD

$$z = H^{1/2} w = \Sigma^{1/2} P w$$

exercise to reader:

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

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

Newton's method for logistic regression: IRLS

... M ... n ...

# Newton's method for logistic regression: IRLS

recall for  $l(w)$ :  $\nabla l(w) = \sum_{i=1}^n x_i [y_i - \sigma(w^T x_i)]$   
 $H(l(w)) = -\sum_{i=1}^n x_i x_i^T \sigma(w^T x_i) \sigma(-w^T x_i)$   
 $V^T H V = -\sum_{i=1}^n \underbrace{(v^T x_i)(x_i^T v)}_{(x_i^T v)^2} \underbrace{\sigma(-) \sigma(-)}_{\neq 0}$

$V^T H V < 0$   $\forall v \in \mathbb{R}^d$   
 i.e. HSO  
 i.e. concave fct.  
 ?  
 Newton is maximizing instead of min.

notation: recall  $X = \begin{pmatrix} -x_1^T \\ \vdots \\ -x_n^T \end{pmatrix}$

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

$\mu_t = \begin{pmatrix} \sigma(w_t^T x_1) \\ \vdots \\ \sigma(w_t^T x_n) \end{pmatrix}$

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

Hessian =  $-\sum_i x_i x_i^T \mu_i (1 - \mu_i) = -X^T D(w) X$  where  $D_{ii} \triangleq \mu_i (1 - \mu_i)$

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

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

$w_{t+1} = (X^T D_t X)^{-1} [X^T D_t z_t]$  where  $z_t \triangleq X w_t + D_t^{-1} (y - \mu_t)$

this is a solution to "weighted least square problem"

$\min_w \| D^{1/2} (z_t - Xw) \|^2$   
 weights      new target

$\sum_i (z_i - w^T x_i)^2$   
 $D_{ii}^{-1}$

compare with Gaussian noise  $\sum_i \frac{(y_i - x_i^T w)^2}{\sigma_i^2}$   
 model for least square

Newton's method for logistic regression

= iterated reweighted least squares (IRLS)

note:  $x = A^{-1}b$   $Ax = b$   $\min_z \|Ax - b\|^2$

# 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 do batch methods  $\nabla f(w) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(w)$   $O(nd)$ 
  - instead you "incremental gradient methods"
  - where  $\nabla f_i(w)$  is gradient of one  $f_i$
  - batch gradient

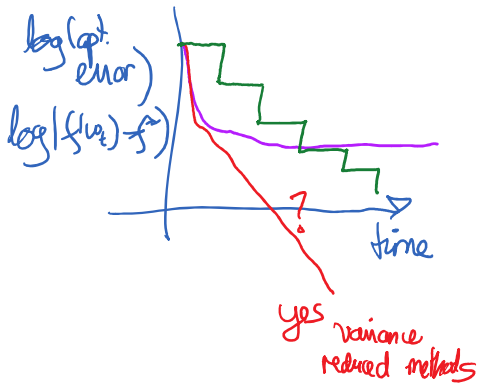
e.g. stochastic gradient descent (SGD) :  $w_{t+1} = w_t - \delta_t \nabla_{S_{i_t}} f(w_t)$   $O(d)$

where  $i_t$  is picked unif. at random

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



batch gradient  $\rightarrow$  expensive updates, but faster convergence



## SAG: stochastic average gradient

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

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

[2012] memory

where  $v_i = \nabla f_i(w_{del})$

at each  $t$ , update  $v_i \triangleq \nabla f_i(w_t)$

$$SAGA: w_{t+1} = w_t - \delta_t \left( \nabla_{S_{i_t}} f(w_t) + \frac{1}{n} \sum_{j=1}^n v_j - v_{i_t} \right)$$

[2014]

variance reduction correction

(default method for log. regression in Scikit-learn)