

today: • logistic regression  
• numerical optimization

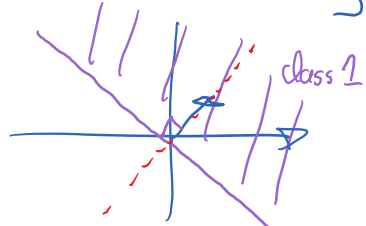
logistic regression model:

$P(y=1|x) = \sigma(w^T x)$   $\mathcal{Y} = \{0, 1\}$  decision rule:  $\mathbb{1}\{w^T x > 0\}$

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

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

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



[if  $\mathcal{Y} = \{\pm 1\}$  "Rademacher" R.V.]

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

encode  $P(y|x) = \sigma(y w^T x)$

max. conditional log-likelihood to estimate  $\hat{w}_{ML}$

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

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

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

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

need to use numerical methods

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

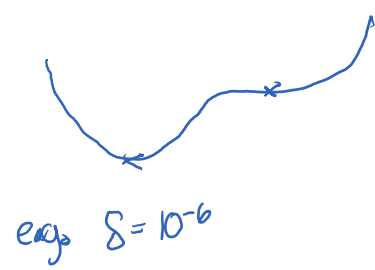
contrast to linear regression  $\nabla l(w) = \sum_{i=1}^n x_i [y_i - w^T x_i]$  linear  $w$

numerical optimization:

want to minimize  $f(w)$  (unconstrained) [suppose  $f$  is diff.]  
 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  $f$   $\mu$ -strongly convex  $\Rightarrow$  if  $\|\nabla f(w_t)\|^2 \leq S$

$\Leftrightarrow f(w) - \frac{\mu}{2} \|w\|^2$  is convex in  $w \Rightarrow f(w_t) - \min_w f(w) \leq \frac{L}{2\mu} \|\nabla f(w_t)\|^2$

step-size rules:

- a) constant step-size:  $\gamma_t = \frac{1}{L}$   $\leftarrow$  Lipschitz continuity constant for  $\nabla f$   
 ie  $\|\nabla f(w) - \nabla f(w')\| \leq L \|w - w'\|$
- b) decreasing step-size rule [this is more common for stochastic optimization]  
 $\gamma_t = \frac{c}{t}$   $\leftarrow$  constant  
 usually want:  $\sum_t \gamma_t = \infty$   
 $\sum_t \gamma_t^2 < \infty$   
 e.g.  $f(w) \triangleq \mathbb{E}_{\xi} g(w, \xi)$   
 $w_{t+1} = w_t - \gamma_t \nabla_w g(w_t, \xi_t)$   
 e.g.  $g(w, \xi) = \ell(y_i, h_w(x_i))$  for ERM  
 $\xi = (x_i, y_i)$

- c) choose  $\gamma_t$  by "line search":  $\min_{\gamma \in \mathbb{R}} f(w_t + \gamma d_t)$   
 $\downarrow$  costly in general  $\uparrow$  direction for update e.g.  $-\nabla f(w_t)$   
 $\rightarrow$  instead do approximate search e.g. Armijo line search (see Boyd's book)

17h20

Newton's method (2nd order method)

motivations minimizing a quadratic approximation

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

Taylor expansion at  $w_t$

$$f(w) = f(w_t) + \nabla f(w_t)^T (w - w_t) + \frac{1}{2} (w - w_t)^T H(w_t) (w - w_t) + \mathcal{O}(\|w - w_t\|^3)$$

$\triangleq Q_t(w)$   $\uparrow$  Taylor's remainder

$= Q_t(w) + \mathcal{O}(\|w - w_t\|^3)$   
 $\uparrow$  quadratic model approx.

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

$$\nabla_w Q_t(w) = 0 \quad \nabla f(w_t) + H(w_t)(w - w_t) \stackrel{\text{want}}{=} 0$$

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

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

$\rightarrow$  inverse Hessian Newton's update  
 $d_t = H^{-1} \nabla f$

~  $O(d^3)$  time to compute in general and  $O(d^2)$  space

$\Leftrightarrow H d_t = \nabla f$

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

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

use numpy.linalg.lstsq  $\uparrow$  to solve this  $\rightarrow$  much more numerically stable

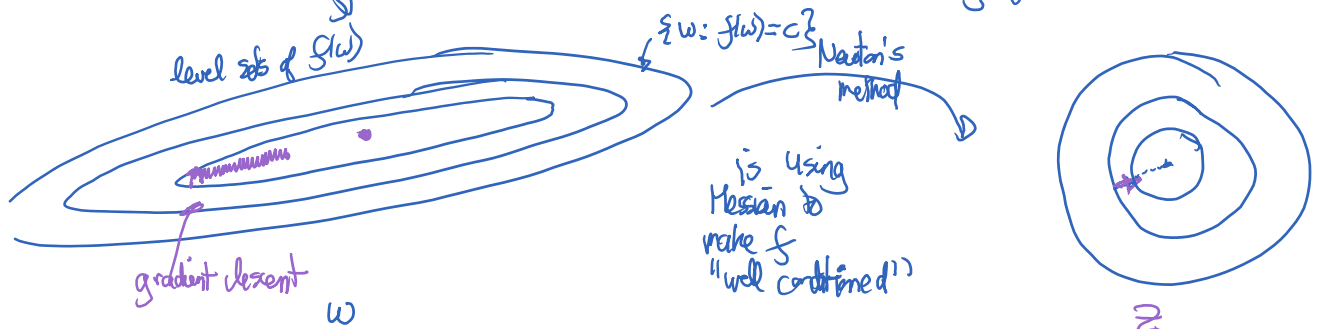
Damped Newton: you add a step-size to stabilize Newton's method

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



Why Newton's method?

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



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

$H = P^T \Sigma P$   
 (H is symmetric)  $\downarrow$  PSD

$\frac{1}{2} w^T P^T \Sigma P w = c \Leftrightarrow \frac{1}{2} z^T z = c$   
diagonal  
 $z = H^{1/2} w = \sum P w$

exercise to read:

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

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

Newton's method for logistic regression: IRLS

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

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

$v^T H v = - \sum_i \underbrace{(v^T x_i)^2}_{\geq 0} \underbrace{\sigma(-) \sigma(-)}_{\geq 0}$

$v^T H v \leq 0 \forall v \neq 0$   
 ie. H  $\leq 0$   
 ie. concave fct

notation: recall  $x_i = \begin{pmatrix} -x_i^T \\ 1 \end{pmatrix}$

Newton's is maximizing instead of min.

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

Newton's is maximizing instead of min.

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

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

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

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 a "weighted least square problem"

$\min_w \left\| \begin{matrix} \sqrt{d_{ii}} \\ \vdots \\ \sqrt{d_{ii}} \end{matrix} (z_t - Xw) \right\|_2^2$

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

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)

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  $\nabla f(w) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(w)$   $O(nd)$   
 (gradient of one datapoint)  
 batch gradient

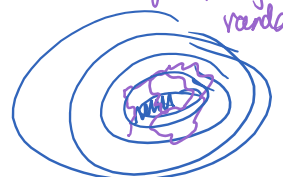
instead you use "incremental gradient methods"

e.g. stochastic gradient descent (SGD) :  $w_{t+1} = w_t - \eta_t \nabla f_{i_t}(w_t)$   $O(d)$

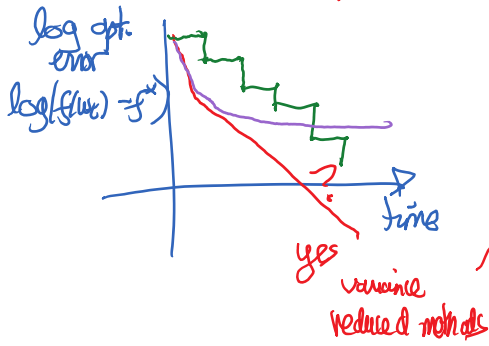
where  $i_t$  is picked unif. at random

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

batch gradient  $\rightarrow$  expensive, but faster convergence



batch gradient  $\rightarrow$  expensive, but faster convergence



SAG: stochastic average gradient

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

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

where  $v_i = \nabla f_i(w_t)$  (word)  
 at each  $t$ , update  $v_i \leftarrow \nabla f_i(w_t)$

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

(default method for logistic regression in Sakit-learn)

variance reduction correction

SVRG