

today: finish Newton + IRLS  
 • Fisher LDA + math tricks

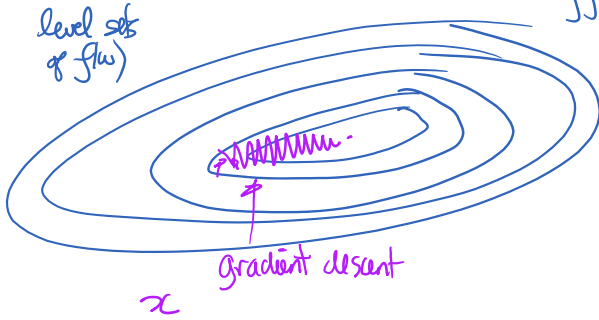
Newton's continuation:

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

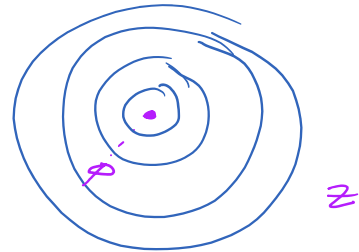
$$w_{t+1} = w_t - \underbrace{\gamma}_{\text{step-size}} H^{-1}(w_t) \nabla f(w_t)$$

why Newton's method?

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



Newton is using Hessian to make  $f$  "well-conditioned"  
 $z = H^{1/2} x$



exercise to readers:

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

$$\downarrow$$

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

"quadratic form"

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

$$\downarrow$$

$$\frac{1}{2} z^T P^T \underbrace{P}_{\text{diagonal}} P x = c$$

$$z = H^{1/2} x = \underbrace{z}_{\substack{P \\ \text{diagonal}}} P x$$

$$\frac{1}{2} z^T z = c$$

$H = P^T \Sigma P$   
 $H$  is symmetric and PSD

Newton's method for logistic regression: IRLS

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

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

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

$v^T H v \leq 0 \quad \forall v \in \mathbb{R}^p$   
 ie.  $H \preceq 0$

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

ie. HJ0  
ie. concave fct.

notation: recall

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

Newton is max.  
instead of min.

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

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

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

↓  
depends on w

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

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

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

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

this is a solution to "weighted least square problem"

$$\min_w \left\| D^{1/2} (z_t - Xw) \right\|^2$$

weights      new target

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

compare with Gaussian noise model for least sq.  $\sum_i (y_i - x_i^T w)^2$

Newton's method for logistic regression

= Iterative reweighted least square (IRLS)

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

Big data logistic regression:

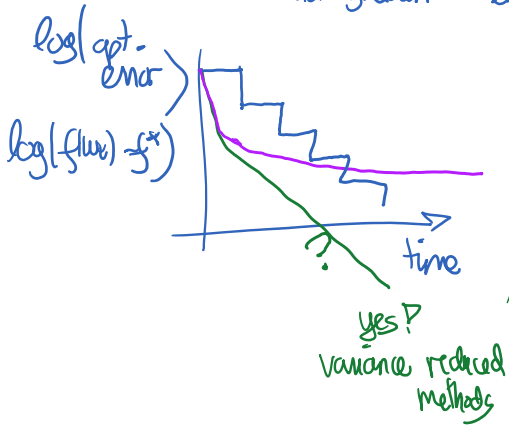
suppose  $f(w) = \frac{1}{n} \sum_{i=1}^n f_i(w)$

- big  $d \Rightarrow$  cannot do  $O(d^2)$  or  $O(d^3)$  operations  $\Rightarrow$  first order methods
- if  $n$  is huge, you cannot do batch methods  $\nabla f(w) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(w)$   $O(nd)$  time  
 instead you "incremental gradient methods" gradient of one fct.  
batch gradient

eg. stochastic gradient descent (SGD):  $w_{t+1} = w_t - \gamma_t \nabla_{w_t} f(w_t)$   $O(d)$  time  
 where  $i_t$  is picked w.r.t  $n$

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

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



SAG: stochastic averaged gradient  
 [2012]

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

SAG:  $w_{t+1} = w_t - \gamma_t \frac{1}{n} \sum_{i=1}^n v_i$  memory  
 where  $v_i = \nabla f_i(w_{t+1})$

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

SAGA:  $w_{t+1} = w_t - \gamma \left( \nabla_{j_t}(w_t) + \frac{1}{n} \sum_j v_j - v_{j_t} \right)$   
 [2014]

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

Variance reduction correction

15/12

generative model for classification: (Fisher) linear discriminant analysis

FLD (LDA)

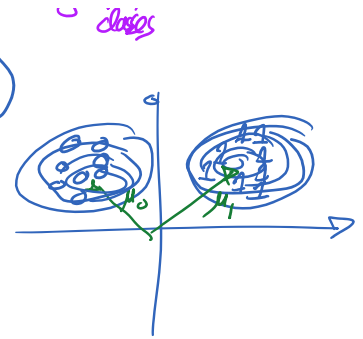
for classification  $Y \in \{0, 1\}$   
 $X \in \mathbb{R}^d$

generative approach:  $p(x, y; \theta) = \overbrace{p(x|y; \theta)}^{\text{class conditionals}} p(y; \theta)$

conditional approach  $p(y|x; \theta)$

for Fisher model: we assume  $p(x|y; \theta) = N(x | \mu_y, \Sigma)$  shared for both classes

for Fisher model: we assume  $p(x|y; \theta) = N(x|\mu_y, \Sigma)$



$$\theta = (\mu_0, \mu_1, \Sigma, \pi)$$

$\uparrow$  mean for class 0       $\downarrow$  shared  $\pi = p(y=1)$

can then show that  $p(y|x; \theta) = \sigma(w^T x)$

where  $w$  is a fct. of  $(\mu_0, \mu_1, \Sigma, \pi)$

[ note: if you use  $\Sigma_0 \neq \Sigma_1$ , get "quadratic discriminant analysis" (QDA) ]  
 i.e.  $\sigma(w^T \phi(x))$  where  $\phi(x)$  is quadratic fct. of  $x$   
 $\rightarrow$  see hwk 2

\* gen. approach: do joint MLE to estimate  $\hat{\theta} = \underset{\theta \in \Theta}{\text{argmax}} \sum_i \log p(x_i, y_i; \theta)$

vs. log. reg. which MLE

side note: MLE for multivariate Gaussian

$X_i \sim N(\mu, \Sigma)$      $\mu \in \mathbb{R}^d$   
 $\Sigma \in \mathbb{R}^{d \times d}$ ,  $\Sigma$  is symmetric,  $\Sigma > 0$

$\Sigma = \mathbb{E}[(x-\mu)(x-\mu)^T]$   
 $\Sigma^T = \mathbb{E}[\text{ " "}] = \Sigma$

$$p(x) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

$\theta = (\mu, \Sigma)$

$\text{tr}(\overbrace{(x-\mu)^T \Sigma^{-1} (x-\mu)}^{\text{inner product}})$   
 $\text{tr}(\Sigma^{-1} (x-\mu)(x-\mu)^T) = \langle \Sigma^{-1}, (x-\mu)(x-\mu)^T \rangle$   
 $\text{tr}(AB) = \text{tr}(BA)$

$\langle A, B \rangle = \sum_{i,j} A_{ij} B_{ij} = \text{tr}(A^T B)$

log-likelihood:  $\sum_{i=1}^n \log p(x_i; \theta) = \text{const.} - \frac{n}{2} \log |\Sigma| - \frac{n}{2} \langle \Sigma^{-1}, \underbrace{\frac{1}{n} \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T}_{\Sigma(\mu)} \rangle$

vector derivative review:

suppose  $f: \mathbb{R}^m \rightarrow \mathbb{R}^n$

$f$  is differentiable at  $x$  iff  $\exists$  linear map  $df: \mathbb{R}^m \rightarrow \mathbb{R}^n$

$\mathcal{O}(\|h\|)$   
 write  $dh$   
 means is

suppose  $f: \mathbb{R}^m \rightarrow \mathbb{R}^n$

$f$  is differentiable at  $x_0$  iff  $\exists$  a linear operator  $df_{x_0}: \mathbb{R}^m \rightarrow \mathbb{R}^n$   
 s.t.  $\forall \Delta \in \mathbb{R}^m$   $f(x_0 + \Delta) - f(x_0) = df_{x_0}(\Delta) + o(\|\Delta\|)$   
↑  
differential

$o(\|\Delta\|)$   
 write as  
 means is  
 some fct.  $h(\|\Delta\|)$   
 s.t.  
 $\lim_{\|\Delta\| \rightarrow 0} \frac{h(\|\Delta\|)}{\|\Delta\|} \rightarrow 0$

$df_{x_0}$  is linear

means  $df_{x_0}(a\Delta_1 + b\Delta_2) = a df_{x_0}(\Delta_1) + b df_{x_0}(\Delta_2)$

can represent as a  $n \times m$  matrix  
 called the Jacobian matrix

standard representation  $(df_{x_0})_{ij} = \frac{\partial f_i}{\partial x_j}$   
↑  
i-th comp. of  $f$   
←  
 $j$ -th component of  $x$

1) this gives a way to get  $df_{x_0}$  for anything (matrix param,  $n$ -dim fct.)

2) be careful with dimensions:  $f: \mathbb{R}^m \rightarrow \mathbb{R}^n$

$df_{x_0}$  is a row vector ( $1 \times m$ )  $df_{x_0} = (\nabla f(x_0))^T$

chain rule: suppose  $f: \mathbb{R}^m \rightarrow \mathbb{R}^n$   
 $g: \mathbb{R}^n \rightarrow \mathbb{R}^q$

$d(g \circ f)_{x_0} = dg_{f(x_0)} \circ df_{x_0}$   
↑  
product of Jacobians

e.g.  $f(x) = x - \mu$   $df_{x_0} = I$

$g(x) = x^T A x$   $dg_{x_0} = x^T (A + A^T)$

$g(f(x)) = (x - \mu)^T A (x - \mu)$   $d(g \circ f)_{x_0} = dg_{f(x_0)} \circ df_{x_0}$   
 $= (x - \mu)^T (A + A^T) \cdot I$

for Gaussian:  $-\frac{1}{2} \sum_i (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)$   
 $\nabla_{\mu}$   $-\frac{1}{2} \sum_i 2 \Sigma^{-1} (x_i - \mu) \stackrel{\text{want}}{=} 0$

$\hat{\mu}_{MLE} = \frac{1}{n} \sum_{i=1}^n x_i$

some references:

- o for convex optimization: [Boyd & Vandenberghe's book](#)

- DL book -- [chapter 4.3 on gradient-based optimization](#)
- for matrix calculus:  
Matrix Differential Calculus with Applications in Statistics and Econometrics, Heinz Neudecker and Jan R. Magnus -- [free online](#)