

today: • brief recap on optimization  
• Fischer model for classification + math tricks

## Numerical optimization

want to minimize  $f(w)$ :  $\min_{w \in \mathbb{R}^d} f(w)$

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

start at  $w_0$

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

step-size rules:

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

b) decreasing step-size rule:  $\gamma_t = \frac{c}{t}$   $\leftarrow$  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  
 $\downarrow$   
(e.g.  $-\nabla f(w_t)$ )

(costly in general,

instead do approximate search  
e.g. Armijo line search

(see Boyd's book) [book link](#)

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

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

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

motivation: minimizing a quadratic approximation:

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

Taylor's remainder

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

quadratic model approximation

want

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

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

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

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

Newton's update

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

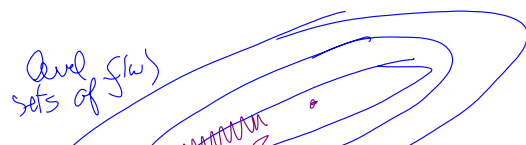
damped Newton: you add a stepsize to stabilize Newton's method

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

why Newton's method?

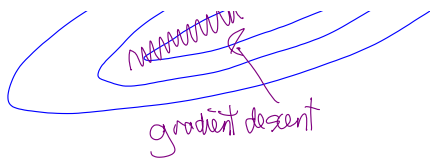
- much faster convergence in # of iterations vs gradient descent
- affine invariant  $\rightarrow$  i.e. invariant to rescaling of variables

but iterations are costly  $O(d^3)$  time  $O(d^2)$  memory

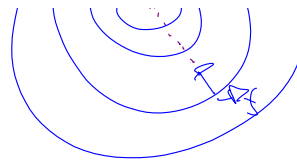


Newton is transforming space using Hessian to make it "more circular"





"well-conditioned"



Newton's method for logistic regression: IRLS

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

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

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

i.e.  $V^T H V \leq 0 \quad \forall V \in \mathbb{R}^d$

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

Let's look at minimizing negative log-likelihood

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

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

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

[note: D depends on w]

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

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

$$\Rightarrow (X^T D_t X)^{-1} [X^T D_t z_t] \quad \text{where } z_t \triangleq X w_t + D_t^{-1} (\vec{y} - \vec{\mu}_t)$$

this is solution to "weighted least square problem"

$$(z_t)_i = x_i^T w_t + \frac{y_i - \mu_i(w_t)}{\mu_i(1 - \mu_i)}$$

$$\min_w \| D^{1/2} (\bar{z}_t - Xw) \|^2$$

new target

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

compare with  
Gaussian noise  
model for least square

$$\sum_{i=1}^n \frac{(y_i - x_i^T w)^2}{\sigma_i^2}$$

data-dependent noise

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

Big data logistic regression

• cannot do  $O(d^2)$  or  $O(d^3)$  operations  $\Rightarrow$  first-order methods

• if  $n$  is huge, you cannot use batch method

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

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

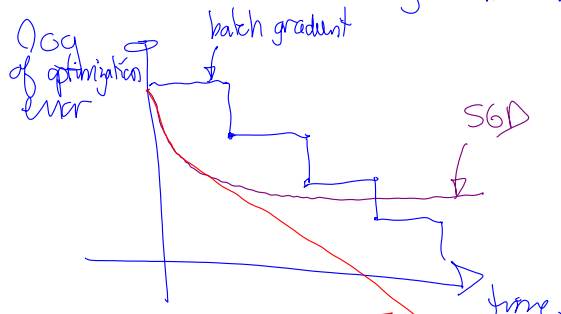
every data point  $\nabla f_i(w)$   $O(nd)$  time  
function depends only on  $i$ th example

instead you use "incremental gradient methods"

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

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

batch gradient  $\rightarrow$  expensive, but fast



SAO: stochastic average gradient  
[2012]

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

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

memory

9 SAB

where  $V_i = \nabla f_i(w_{old})$

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

SAGA :  $w_t \leftarrow \gamma (\nabla f_{i_t}(w_t) + \underbrace{\frac{1}{n} \sum_{i=1}^n V_i - V_{i_t}}_{\text{variance reducing correction}})$

[2014] [paper](#)

(default method for logistic regression in Scikit-learn)

generative model for classification: (Fisher) linear discriminant analysis

for classification  $y \in \{0, 1\}$   
 $x \in \mathbb{R}^d$

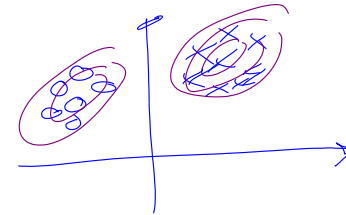
generative approach  $p(x, y; \theta) = \overbrace{p(x|y; \theta)}^{\text{class conditions}} p(y; \theta)$

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

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

same for both classes

$\mu_0, \mu_1$



$\theta = (\mu_0, \mu_1, \underbrace{\Sigma}_{d \times d \text{ shared for 2 classes}}, \underbrace{\pi_0}_{\text{prior } p(y=1)})$

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

where  $w$  is function of  $(\mu_0, \mu_1, \Sigma, \pi)$

[if you use  $\Sigma_0$  &  $\Sigma_1$ , get "quadratic discriminant analysis" (QDA) → see book. 2]

[if you use  $\Sigma_0$  &  $\Sigma_1$ , get "quadratic discriminant analysis" (QDA)  $\rightarrow$  see book 2]

i.e.  $\sigma(w^T \phi(x))$  where  $\phi(x)$  is quadratic in  $x$

⊗ generative approach: do joint ML i.e.  $\hat{\Theta} = \arg \max_{\Theta \in \Theta} \sum_i \log p(x_i, y_i; \Theta)$

side note: MLE for multivariate Gaussian

$$x_i \sim N(\mu, \Sigma) \quad \mu \in \mathbb{R}^d$$

$$\Sigma \in \mathbb{R}^{d \times d}, \quad \Sigma \text{ is symmetric} \\ \Sigma \succ 0$$

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

$$\text{tr}(AB) = \text{tr}(BA)$$

$$\begin{aligned} & \text{tr}((x-\mu)^T \Sigma^{-1} (x-\mu)) \\ &= \text{tr}\left(\Sigma^{-1} (x-\mu)(x-\mu)^T\right) = \left\langle \Sigma^{-1}, (x-\mu)(x-\mu)^T \right\rangle \\ & \langle A, B \rangle \triangleq \sum_{i,j} A_{ij} B_{ij} = \text{tr}(A^T B) \end{aligned}$$

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

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

vector derivative review:

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.  $\forall \Delta \in \mathbb{R}^m, \quad f(x_0 + \Delta) - f(x_0) = df_{x_0}(\Delta) + o(\|\Delta\|)$

"little oh"

$o(\|\Delta\|)$  means some function  $h(\|\Delta\|)$  s.t.

$$\lim_{\|\Delta\| \rightarrow 0} \frac{h(\|\Delta\|)}{\|\Delta\|} \rightarrow 0$$

$$f(x) \in \mathbb{R}^n$$

$$\Delta \in \mathbb{R}^m$$

$$\text{i.e. } df_{x_0}: \mathbb{R}^m \rightarrow \mathbb{R}^n$$

$$\text{is linear i.e. } df_{x_0}(\Delta_1 + b\Delta_2) = df_{x_0}(\Delta_1) + b df_{x_0}(\Delta_2)$$

can represent as a  $n \times m$  matrix

called the Jacobian matrix

$$\text{standard representation } (df_{x_0})_{ij} = \frac{\partial f_i}{\partial x_j} \quad \begin{matrix} \swarrow \text{i-th component of } f \\ \nwarrow \text{j-th component of } x \end{matrix}$$

1) This gives you a way to get  $df_{x_0}$  for anything (matrix, tensor, ...)

2) careful with dimension: differential of  $f: \mathbb{R}^m \rightarrow \mathbb{R}$  is a row vector ( $1 \times m$ ) i.e.  $df_{x_0} = (\nabla f(x_0))^T$

chain rule:

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

$$g: \mathbb{R}^n \rightarrow \mathbb{R}^q$$

$$d(g \circ f)_{x_0} = dg_{f(x_0)} \circ df_{x_0}$$

e.g.

$$f(x) = x - \mu \quad df_{x_0} = I$$

$$g(x) = x^T A x \quad dg_{x_0} = x^T (A + A^T) \quad [\nabla g^T]$$

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

for Gaussian

$$-\frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)$$

$\nabla_{\mu}$

$$-\frac{1}{2} \sum_{i=1}^n 2 \Sigma^{-1} (x_i - \mu) \stackrel{\text{want}}{=} 0$$

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

2.5.1

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

example 2: derivative of  $f(A) \triangleq \log \det(A)$

can represent derivative of a fun. from matrix to scalar, as a matrix

$$\begin{aligned} f(A+\Delta) - f(A) &= \text{tr}(f'(A)^T \Delta) + o(\|\Delta\|) \\ &= \langle f'(A), \Delta \rangle + o(\|\Delta\|) \end{aligned}$$

$$\log \det(A+\Delta) - \log \det A \quad \left[ \text{assume } A \text{ is symmetric} \right]$$

$$= \log \det(A^{1/2} (I + A^{-1/2} \Delta A^{-1/2}) A^{1/2}) - \log \det A$$

$$= \log \cancel{|A|^{1/2}} |I + A^{-1/2} \Delta A^{-1/2}| \cancel{|A|^{1/2}} - \log |A|$$

$$= \log \det(I + A^{-1/2} \Delta A^{-1/2})$$

$$= \sum_i \log(1 + \lambda_i(A^{-1/2} \Delta A^{-1/2}))$$

$$\leftarrow \det(A) = \prod_i \lambda_i(A) \quad \begin{array}{c} \text{e-values} \\ \downarrow \end{array}$$

$$= \sum_i \lambda_i(A^{-1/2} \Delta A^{-1/2}) + o(\|\Delta\|)$$

$$\leftarrow \log(1+x) = x + o(x^2) \text{ for } |x| < 1$$

$$= \text{tr}(A^{-1/2} \Delta A^{-1/2}) + o(\|\Delta\|)$$

$$\leftarrow \text{tr}(A) = \sum_i \lambda_i(A)$$

$$= \text{tr}(A^{-1} \Delta) + o(\|\Delta\|)$$

$$\Rightarrow \boxed{\frac{d}{dA} \log \det(A) = A^{-1}}$$

back to our log-likelihood:



derivative  
w.r. to  $\Sigma^{-1}$

$$\frac{n}{2} \log |\Sigma^{-1}| - \frac{n}{2} \langle \Sigma^{-1}, \tilde{\Sigma} \rangle$$

$$\frac{n}{2} \underbrace{(\Sigma^{-1})^{-1}}_{\tilde{\Sigma}} - \frac{n}{2} \tilde{\Sigma} \stackrel{\text{const}}{=} 0 \Rightarrow \hat{\Sigma}_{\text{MLE}} = \tilde{\Sigma} = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_{\text{MLE}})(x_i - \mu_{\text{MLE}})^T$$

**some references:**

- 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](#)