

(1)

Solutions to nonlinear equations

Let $f: \mathbb{R} \rightarrow \mathbb{R}$. Solve $f(x) = 0$ for $x \in [a, b]$

- Bisection

functional iteration

- Newton's method

Stat Model

Data

Technique
+ principle

Statistics

Algorithm

Program

(1) Probability

(2) Linear Algebra

(3) Optimization ←

We look at (3).

Generally, we want to maximize or minimize something.

 \Rightarrow max! likelihood \Rightarrow min! sum of squares \Rightarrow ~~max f~~ \Rightarrow $\max f = \min -f$
So don't worry about it.

This course is about

- (1) maximizing a function
- (2) integrating a function

Course outline ① Solving non linear equations (root finding)

$$f(x) = 0 \quad \text{for } f: \mathbb{R} \rightarrow \mathbb{R}$$

$$f(x) = 0 \quad \text{for } f: \mathbb{R}^K \rightarrow \mathbb{R}$$

$$k=2, 3, \dots, N \quad (+\text{other})$$

② General optimization routines

$$\text{Given } f: \mathbb{R}^K \rightarrow \mathbb{R}, \quad \max_x f(x)$$

$$\text{or } \min_x f(x)$$

Line search Methods

↳ Newton

↳ Quasi-Newton

Taylor's theorem

① Pick a direction

② Go a certain distance M that direction

related: Simulate annealing

A random optimizer — still general purpose

③ Statistics!

EM algorithm for maximum likelihood

Minimization / Majorization

Monte Carlo EM

~~Deterministic Algorithms~~

General idea:

Entire
class/course

- ① It's difficult to ~~optimize~~ ^{maximize} f .
 - ② We can compute an approximation to f called g .
 - ③ "Transfer" optimization to g and ~~maximize~~ ^{maximize} g .
 - ④ Iterate ②, ③
- ⇒ Instead of direct max, "transfer" to simpler function and iterate



$$f(b) = f(a) + f'(a)(b-a)$$

$$f(x_n) = f(x_0) + f'(x_0)(x_n - x_0)$$

$$f(x_{n+1}) = f'(c)(x_{n+1} - x_0)$$

~~some outline~~

Course outline

Integration
 "How to compute
 an integral"

density $\Rightarrow p(x) = C f(x)$
 ↗
 integration

Analytic approximation - Laplace approx.
 Quadrature

* Monte Carlo (integration)

(EM algorithm = "avoiding integrals")
 Random numbers, rejection sampling
 Importance sampling

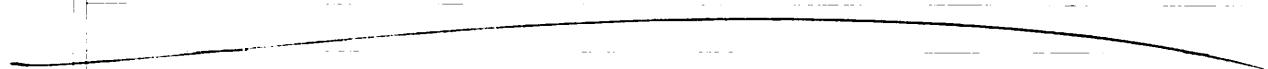
* Markov chain Monte Carlo (MCMC)

draw samples from a posterior distribution
 Metropolis-Hastings
 Gibbs sampling
 Variants / ~~tricks~~

* Smoothing

↳ Splines, Kernel smoothing, P-splines
 linear smoothers
 (gams)

* Bootstrap



Miscellaneous

Computation

① Linear regression

R uses QR decom

$$\hat{\beta} = (X'X)^{-1} X' y \quad X = QR$$

$$\hat{\beta} = \cancel{(R'R)^{-1} R'Qy}$$

②

$$X'X\beta = X'y$$

② Multivariate normal

$$R'Q'Q^*R\beta = R'Q'y$$

$$R'R\beta = R'Q'y$$

$$q(x|\mu, \Sigma)$$

$$R\beta = Q'y$$

$$= -\frac{1}{2} \underbrace{\log |\Sigma|}_{\det} - \frac{1}{2} (x-\mu)' \Sigma^{-1} (x-\mu)$$

$$q(z|\Sigma) = -\frac{p}{2} \log 2\pi - \log |\Sigma| - \frac{1}{2} z' \Sigma^{-1} z$$

$$\cancel{\Sigma} = R'R \quad (\text{cholesky})$$

$$z' \Sigma^{-1} z = z'(R'R)^{-1} z : z' R^{-1} R'^{-1} z$$

$$\begin{aligned}
 &= z' \underbrace{R'^{-1}}_x \underbrace{R'^{-1}}_x z \\
 &= (R'^{-1} z)' R'^{-1} z \\
 R \cancel{x} &= \cancel{x} \quad \text{backsolve}
 \end{aligned}$$

$$X = \underbrace{R'^{-1}}_{R'^{-1} R R'^{-1} R'^{-1}} z$$

$$\Rightarrow z' \Sigma^{-1} z = x' x$$

$$\begin{aligned}
 &z' (R'R)^{-1} z \\
 &z' R^{-1} R'^{-1} z \\
 &= (R'^{-1} z)' R'^{-1} z \\
 X' &= z' R
 \end{aligned}$$

Solving nonlinear equations

$$f(x) = 0 \quad \text{for } x \in [a, b]$$

Bisection method.

~~Method of false position~~

If $\text{sign}(f(a)) \neq \text{sign}(f(b))$

\Rightarrow Intermediate value theorem

Let $f(a) < \gamma < f(b)$. $\exists c \in [a, b]$
s.t. $f(c) = \gamma$.

\Rightarrow i.e. If f is cont. on $[a, b]$, $f^{-1}([a, b])$ is closed

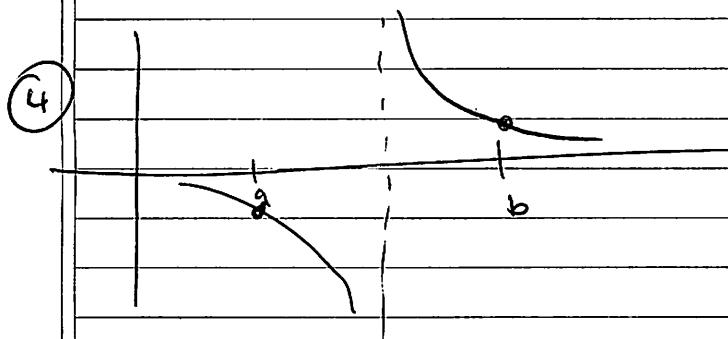
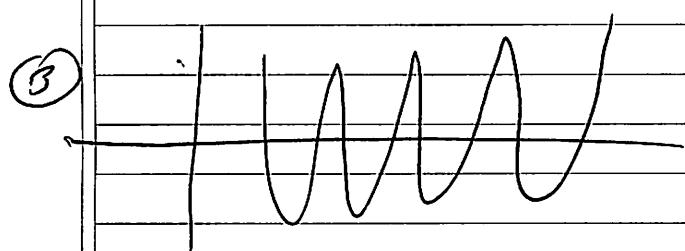
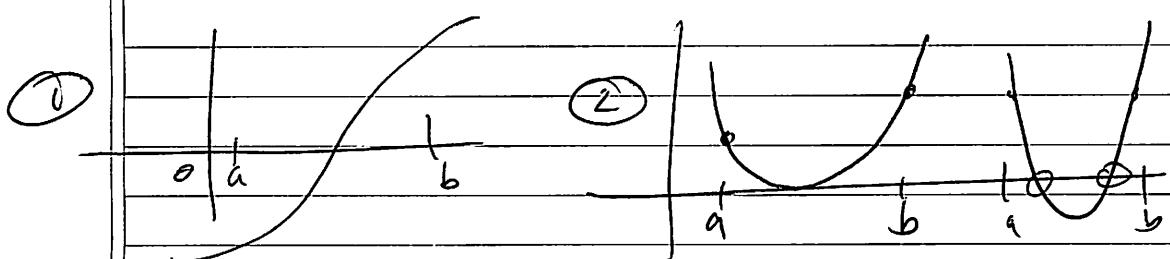
① Let $c = \frac{a+b}{2}$

② If $f(c) = 0$, stop

③ Else if $\text{sign}(f(a)) \neq \text{sign}(f(c))$, $b \leftarrow c$.
else if $\text{sign}(f(b)) \neq \text{sign}(f(c))$, $a \leftarrow c$.

④ goto ①

For n iterations, size of interval $\approx 2^{-n}(b-a)$



Converge when $|b-a| < \varepsilon$ or
 $|f(b) - f(a)| < \varepsilon$. Depends on situation.

Ex.

$$l(\theta) = \text{likelihood}$$

$$l'(\theta) = 0 \Rightarrow \hat{\theta} \text{ is MLE}$$

Often want $|l(b) - l(a)| < \varepsilon$ even if l is flat.

Ex. Quantiles

Given Cdf $F(x)$, ~~want to find x s.t.~~

and prob $p \in (0, 1)$, find x s.t. $F(x) = p$.

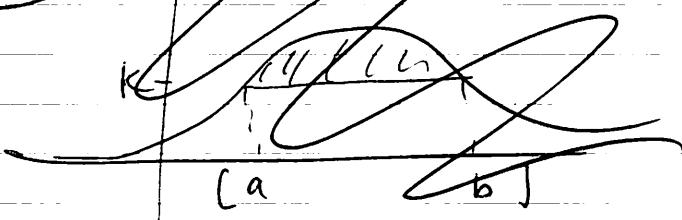
Let $g(x) = F(x) - p$.

Solve $g(x) = 0$

~~Ex. likelihood intervals~~ ~~HPD intervals~~
Bayesian Credible Intervals

Given K ,

$$S_K \stackrel{\Delta}{=} \{ \theta : f(\theta | x) \geq K \}$$



Ex Likelihood Intervals

~~Let~~ Let $f(\theta) = L(\theta)/L(\hat{\theta})$

Find $\text{LI} = \{ \theta : f(\theta) \geq \gamma_8 \}$

Solve $f(\theta) - \gamma_8 = 0$

Ex. Bayesian credible interval

$$\text{Let } S_K = \{\theta : f(\theta | y) \geq K\}$$

Bayesian credible interval of level $\alpha \Leftrightarrow \text{s.t. } S_K$

~~$R(\theta \in S_K | y) = \alpha$~~

$$\mu([a, b]_K) = \alpha$$

$$\mu([a, b]_K) - \alpha = 0 \quad \text{Solve for } K$$

$$[a, b]_K$$

Solve for a, b .

For $f: \mathbb{R}^K \rightarrow \mathbb{R}, K=2, 3, \dots, N$

Initial "box" area = $\prod_{i=1}^K (b_i - a_i)$

At iteration n , area $\approx \frac{1}{2} (b_i - a_i)$

~~$\text{area} = \frac{1}{2^n} \prod_{i=1}^K (b_i - a_i)$~~

Correction algorithm: ~~as interval length of $\frac{1}{2^n}$~~

At iteration 1, area = $\prod_{i=1}^K \frac{1}{2} (b_i - a_i) = \frac{1}{2^K} \prod_{i=1}^K (b_i - a_i)$

2: area = $\prod_{i=1}^K \frac{1}{2^2} \frac{1}{2} (b_i - a_i) = \frac{1}{2^K} \prod_{i=1}^K (b_i - a_i)$

\vdots area = $\frac{1}{2^n} \prod_{i=1}^K (b_i - a_i)$

Rates of convergence

(1)

Suppose $X_n \rightarrow X_\infty$ in \mathbb{R}^k . ~~Q-linear~~

Say the convergence is Q-linear ("linear") if $\exists r \in (0, 1)$

$$\frac{\|X_{n+1} - X_\infty\|}{\|X_n - X_\infty\|} \leq r \text{ for all } n \text{ sufficiently large.}$$

Ex. ~~$X_n = 1 + 2^{-n}$~~ $X_n = 1 + 2^{-n}$ is Q-linear.
 $X_\infty = 1$

(2)

~~Q-quadratic~~
Q-super linear if

$$\lim_{n \rightarrow \infty} \frac{\|X_{n+1} - X_\infty\|}{\|X_n - X_\infty\|} = 0$$

Ex. $X_n = 1 + n^{-1}$ is Q-super linear

(3)

Q-Quadratic if

$$\frac{\|X_{n+1} - X_\infty\|}{\|X_n - X_\infty\|^2} \leq M \text{ for all } n \text{ suff. large}$$

Ex. $X_n = 1 + 2^{-2^n}$

Ex. M bisection algorithm:

let $x_n = |b_n - a_n|$, i.e. size of interval at iteration n. Then

$$\frac{|x_{n+1} - x_\infty|}{|x_n - x_\infty|} = \frac{x_{n+1}}{x_n} = \frac{2^{-(n+1)}(b_0 - a_0)}{2^{-n}(b_0 - a_0)} \\ = 2^{-n-1+n} = \frac{1}{2} \leq r \in (0,1)$$

Bisection achieves linear convergence

~~Quasi Newton~~

Newton's Method — quadratic

Quasi-Newton — superlinear

steepest descent — linear

~~Newton~~ \Rightarrow

Functional Iteration

We want to solve $f(x) = 0$ for $f: \mathbb{R}^K \rightarrow \mathbb{R}$
and $x \in S \subset \mathbb{R}^K$

Any root of f is a fixed point of

$g(x) = f(x) + x$. (There are other functions)

$g(x) = x(f(x) + 1)$, $x \neq 0$

Solutions to $f(x) = 0$ are fixed points of
other functions.

Sometimes we can take a function f and create a sequence $X_n = f(X_{n-1})$. Depending on f , we can have $X_n \rightarrow X_\infty$ where $f(X_\infty) = X_\infty$ (i.e. a fixed point).

When does iteration functional work?

$\langle\langle$ Shrinking lemma $\rangle\rangle$

Newton's Method

Solve $f(x) = 0$. Get solution X_∞ and let X_n be our current estimate. By MVT,

$$f(X_n) = f'(z)(X_n - X_\infty) \text{ where}$$

$z \in$ b/w X_n and X_∞ .

$$\Rightarrow X_\infty = X_n - \frac{f(X_n)}{f'(z)}$$

Since X_∞ and z unknown, ~~let's do~~

$$X_{n+1} = X_n - \frac{f(X_n)}{f'(X_n)}$$



Newton update

$\langle\langle$ Proof of Newton's Method $\rangle\rangle$

Shrinking Lemma

$$s = 1 + \frac{1}{r} + \frac{1}{r^2} + \dots$$

$$rs = r + \frac{r}{r} + \frac{r}{r^2} + \dots$$

$$rs = r + s$$

$$\cancel{s(r-1) \in S}$$

$$r \cdot \cancel{s} \in S$$

$$s(r-1) = r$$

$$s = \frac{r}{r-1} = \frac{1}{1-\frac{1}{r}}$$

Let M be a closed subset of a c.n.v.s. let $f: M \rightarrow M$ be a map, and assume $\exists K, 0 < K < 1$ s.t. $\forall x, y \in M$, we have

$$|f(x) - f(y)| \leq K|x - y|.$$

Then f has a unique fixed point. There is a unique pt. $x_0 \in M$ s.t. $f(x_0) = x_0$.

\Rightarrow If $x \in M$, the sequence $\{f^n(x)\}$ is a Cauchy sequence which converges to x_0 .

Prf:

Given $x \in M$, we have

$$|f^2(x) - f(x)| = |f(f(x)) - f(x)| \leq K|f(x) - x|.$$

By induction:

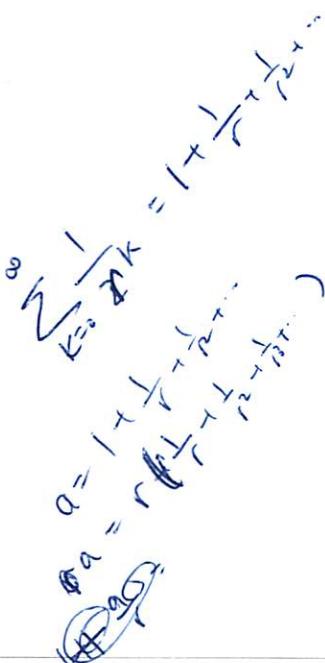
$$|f^{n+1}(x) - f^n(x)| \leq K|f^n(x) - f^{n-1}(x)| \leq K^n|f(x) - x|$$

And the set of elements $\{f^n(x)\}$ is bounded because

$$|f^n(x) - x| \leq |f^n(x) - f^{n-1}(x)| + |f^{n-1}(x) - f^{n-2}(x)| + \dots + |f(x) - x|$$

$$\leq \underbrace{(K^{n-1} + K^{n-2} + \dots + K)}_{\text{geometric series}} |f(x) - x|$$

$$\leq \frac{1}{1-K} |f(x) - x|$$



$$\left| f^{m+k}(x) - f^m(x) \right| \leq$$

f^m

By induction, Given $m \geq 1, k \geq 1$, we have

$$\left| f^{m+k}(x) - f^m(x) \right| \leq K^m \underbrace{\left| f^k(x) - x \right|}_{\leq \frac{1}{1-k} |f(x) - x|}$$

$\Rightarrow \exists N$ s.t. if $m, n \geq N$ (say $n = m+k$),

$$\left| f^{m+k}(x) - f^m(x) \right| < \varepsilon$$

because $K^m \rightarrow 0$ as $m \rightarrow \infty$.

$\Rightarrow \{f^n(x)\}$ is a Cauchy sequence. Let x_0 ~~be its~~ ^{limit} be its limit. Let N be s.t. $\forall n \geq N, |f^n(x) - x_0| < \varepsilon$.

Then

$$\left| f(x_0) - f^{n+1}(x_0) \right| \leq K \left| x_0 - f^n(x_0) \right| < \varepsilon$$

$$\Rightarrow \{f^n(x)\} \rightarrow f(x_0), \{f^n(x)\} \rightarrow x_0$$

$\Rightarrow f(x_0) = x_0$, a fixed point

Let x_1 be another fixed point. Then

$$|x_1 - x_0| = |f(x_1) - f(x_0)| \leq K |x_1 - x_0|.$$

Since $0 < K < 1$, $x_1 = x_0$, hence unique \square

Thm:

Let $f \in C^2$ and suppose $\exists x_0$ s.t. $f(x_0) = 0$ and $f'(x_0) \neq 0$. Then $\exists \delta$ s.t. for any $x_0 \in [x_0 - \delta, x_0 + \delta]$, & the sequence

$$x_n = g(x_{n-1}) = x_{n-1} - \frac{f(x_{n-1})}{f'(x_{n-1})}$$

converges to x_0 !

Pmf:

Note that

$$\begin{aligned} g'(x) &= \cancel{\left(-\frac{f(x)f''(x) - f'(x)f'(x)}{[f'(x)]^2} \right)} \quad \cancel{+ \frac{f'(x)f'(x) - f(x)f''(x)}{[f'(x)]^2}} \\ &= \cancel{\frac{f(x)f''(x)}{[f'(x)]^2}} \end{aligned}$$

$$\Rightarrow g'(x_0) = 0$$

Since $f \in C^2$, g' is continuous. Therefore g is

Given $k < 1$, $\exists \delta > 0$, s.t. $\forall x \in [x_0 - \delta, x_0 + \delta] = A$

$$|g'(x)| < k.$$

Also, Given $a, b \in A$,

$$\begin{aligned} |g(a) - g(b)| &\leq |g'(c)| |a - b| \\ &\leq k |a - b| \quad (0 < k < 1) \end{aligned}$$

$\Rightarrow g$ is a shrinking map. on A .

$\Rightarrow \exists$ unique x_0 s.t. $g(x_0) = x_0$

$\max f = \text{classical/frequentist}$
 $\int f d\mu = \text{Bayesian}$

Convergence rates for shrinking maps.

~~Suppose $g: \mathbb{R}^k \rightarrow \mathbb{R}$ and~~

Suppose g satisfies

$$|g(x) - g(y)| \leq K|x - y|$$

for some $K \in (0, 1)$ and any $x, y \in I$, a closed interval.

Also, assume $0 < |g'(x_0)| < 1$, where

x_0 is the fixed point. Then $x_n \rightarrow x_0$ at a linear rate.

$$\text{PF: } \frac{|x_{n+1} - x_0|}{|x_n - x_0|} = \frac{|g(x_n) - g(x_0)|}{|x_n - x_0|} \leq K \frac{|x_n - x_0|}{|x_n - x_0|} = K$$

$\Rightarrow K$ is a limit

$$\lim_{n \rightarrow \infty} \frac{|g(x_n) - g(x_0)|}{|x_n - x_0|} = |g'(x_0)| \geq 0$$

constant $\in (0, 1)$

\Rightarrow linear convergence

What about Newton's method?

Suppose $f \in C^2$, ~~and~~ and $\exists x_\infty$ s.t.
 $f(x_\infty) = 0$.

By Taylor's theorem: for some small ε ,

$$\textcircled{1} \quad f(x_\infty + \varepsilon) = f(x_\infty) + \varepsilon f'(x_\infty) + \frac{\varepsilon^2}{2} f''(x_\infty) + O(\varepsilon^3)$$

~~$$\textcircled{2} \quad f(x_\infty + \varepsilon) = 0 + \varepsilon f'(x_\infty) + \frac{\varepsilon^2}{2} f''(x_\infty) + O(\varepsilon^3)$$~~

$$\textcircled{2} \quad f'(x_\infty + \varepsilon) = f'(x_\infty) + \varepsilon f''(x_\infty) + O(\varepsilon)$$

Newton's method generates the sequence

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$

$$\Rightarrow x_{n+1} - x_\infty = x_n - x_\infty - \frac{f(x_n)}{f'(x_n)}$$

$$\text{Let } \varepsilon_{n+1} = x_{n+1} - x_\infty, \quad \varepsilon_n = x_n - x_\infty.$$

$$\Rightarrow \varepsilon_{n+1} = \varepsilon_n - \frac{f(x_n)}{f'(x_n)}$$

$$\text{By +/-, } \varepsilon_{n+1} = \varepsilon_n - \frac{f(x_\infty + \varepsilon_n)}{f'(x_\infty + \varepsilon_n)}$$

$$\Rightarrow \varepsilon_{n+1} = \varepsilon_n - \frac{\varepsilon_n f'(x_\infty) + \varepsilon_n^2 f''(x_\infty)/2}{f'(x_\infty) + \varepsilon_n f''(x_\infty)}$$

$$= \frac{\cancel{\varepsilon_n f'} + \cancel{\varepsilon_n^2 f''} - \cancel{\varepsilon_n f'} - \cancel{\varepsilon_n^2 f''}/2}{\cancel{f'} + \varepsilon_n f''}$$

$$= \varepsilon_n^2 \left(\frac{f''/2}{f' + \varepsilon_n f''} \right)$$

$$\frac{\varepsilon_{n+1}}{\varepsilon_n^2} \approx \frac{f''(x_\infty)/2}{f'(x_\infty) + \varepsilon_n f''(x_\infty)}$$

$$\approx \frac{f''(x_\infty)/2}{f'(x_\infty)}$$

$\varepsilon_n \downarrow 0$

$\Rightarrow \exists$ some $M < \infty$ s.t.

$$\left| \frac{\varepsilon_{n+1}}{\varepsilon_n^2} \right| \leq M, \quad \forall n \text{ suff. large.}$$

\Rightarrow Quadratic convergence.

Of course we need $f''(x_\infty)$ exists and $f'(x_\infty) \neq 0$.

In practice, we ignore assumptions / conditions.
Use Newton's method as a "black box". Caution Empirical.

Pro : Very fast in neighbourhood of truth
Direct multivariate generalization

Con : Need to evaluate f'
 $\&$ can be unstable.

We want $\hat{\theta}$, the value of θ that maximizes $\ell(\theta)$. Assume that $\hat{\theta}$ is the unique root of $\ell'(\theta) = 0$. Solve $\ell'(\theta) = 0$ (likelihood equations)

Newton's method

$$\theta_{n+1} = \theta_n - [\ell''(\theta_n)]^{-1} \ell'(\theta_n)$$

$K \times 1$ $K \times K$ $K \times 1$

\Rightarrow May be easier/better to solve

$$[\ell''(\theta_n)]\theta_{n+1} = [\ell''(\theta_n)]\theta_n - \ell'(\theta_n)$$

$A \quad x = b$

Then try to invert $\ell''(\theta_n)$.

At convergence we have, in addition to $\hat{\theta}$,

$\ell'(\hat{\theta})$: score statistic

$-\ell''(\hat{\theta})$: observed information

The obs. information is related to the ~~covariance~~ matrix of ~~uniting~~ normal dist. of $\hat{\theta}$, i.e.

$$\sqrt{n} \left([-\ell''(\hat{\theta})]^{1/2} (\hat{\theta} - \theta_0) \right) \rightarrow N(0, I)$$

for $n \rightarrow \infty$.

Binary logistic regression

$$\log p_i - \log(1-p_i) = x_i^\top \beta$$

$$\log p_i = x_i^\top \beta$$

$$Y_i \sim \text{Bernoulli}(p_i), i=1, \dots, n$$

$$\text{logit}(p_i) = \log \frac{p_i}{1-p_i} = x_i^\top \beta$$

$$p_i = \frac{e^{x_i^\top \beta}}{1+e^{x_i^\top \beta}}$$

$$L(\beta) \propto \prod_{i=1}^n p_i^{Y_i} (1-p_i)^{1-Y_i}$$

$$= \exp\left(\sum_{i=1}^n Y_i \log p_i + (1-Y_i) \log(1-p_i)\right)$$

$$= \cancel{\exp\left(\sum_{i=1}^n Y_i \log(x_i^\top \beta) + (1-Y_i) \log(1+e^{x_i^\top \beta})\right)}$$

$$= \exp\left(\sum Y_i (x_i^\top \beta - \log(1+e^{x_i^\top \beta})) + (1-Y_i) (-\log(1+e^{x_i^\top \beta}))\right)$$

$$\frac{\partial L}{\partial \beta} = \sum Y_i \left[x_i - \frac{x_i e^{x_i^\top \beta}}{1+e^{x_i^\top \beta}} \right] + (1-Y_i) \left[\frac{-x_i e^{x_i^\top \beta}}{1+e^{x_i^\top \beta}} \right]$$

$$= \sum Y_i [x_i - x_i p_i] + (1-Y_i) (-x_i p_i) \\ = \sum Y_i [x_i - x_i p_i + x_i p_i] - x_i p_i \\ = \sum Y_i [x_i - x_i p_i] - x_i p_i$$

$$l''(\beta) = -X^T W X$$

$$W = \text{diag}(p_i(1-p_i))$$

$$= \sum Y_i x_i - x_i p_i$$

$$= \sum x_i^T (Y_i - p_i)$$

$$= \cancel{\sum x_i^T (p_i(1-p_i) - p_i^2)}$$

$$= X^T (Y - p)$$

TOPIC:

DATE:

FILE UNDER:

PAGE:

$$e'(\beta) = X^T(Y - P)$$

$$e''(\beta) = -X^T W X$$

$$\hookrightarrow \text{diag}[P_i(1-P_i)]$$

$$\beta_{n+1} = \beta_n + [-X^T W X]^{-1} [X^T(Y - P_n)]$$

Newton update

For exp. families
w/canonical link
 $\hat{E}l''(\theta) = E l''(0)$
so Newton and
Fisher scoring are same

General Purpose Minimization

Given a function $f: \mathbb{R}^k \rightarrow \mathbb{R}$, we want to
find $\min_{x \in S} f(x)$ where $S \subset \mathbb{R}^k$.
usually $\Leftrightarrow f'(x) = 0$

Line search methods

Given f and a current estimate of the location
of the minimum x_n , we want to

① Choose a direction p_n (vector)

② Solve $\min_{\alpha > 0} f(x_n + \alpha p_n)$

↳ don't need exact min. Rather compute some
candidates and choose the best one.

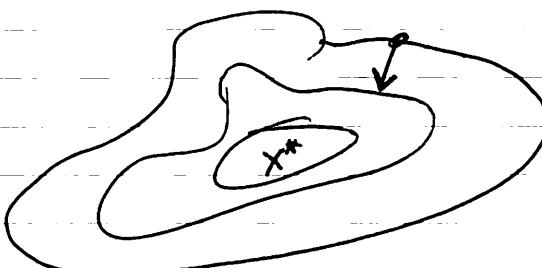
③ $x_{n+1} = x_n + \alpha p_n$

Choosing Direction

— Most obvious \Rightarrow steepest-descent : $\nabla f(x_n)$

$-f'(x_n)$ direction along which f decreases most
rapidly

↳ orthogonal to contours of f



Newton direction:

By Taylor's Theorem:

$$f(x_n + p) \approx f(x_n) + p^T f'(x_n) + \frac{1}{2} p^T f''(x_n) p$$

↓
 $m_n(p)$

Minimize $m_n(p) \rightarrow p_n = [-f''(x_n)]^{-1} f'(x_n)$
over p

{ Newton direction has "natural" step length of 1 }
but this can be modified

$$\Rightarrow x_{n+1} = x_n + [-f''(x_n)^{-1}] f'(x_n)$$

That's familiar!

Similarly, Quasi-Newton

$$p_n = B_n^{-1} f'(x_n)$$

$$f'(x_n) - f'(x_{n-1}) = B_n(x_n - x_{n-1})$$

B_n satisfies a "secant condition"

Coordinate descent

If f is K -dimensional, we ~~can't~~ minimize along M dimensions in a cyclic fashion.

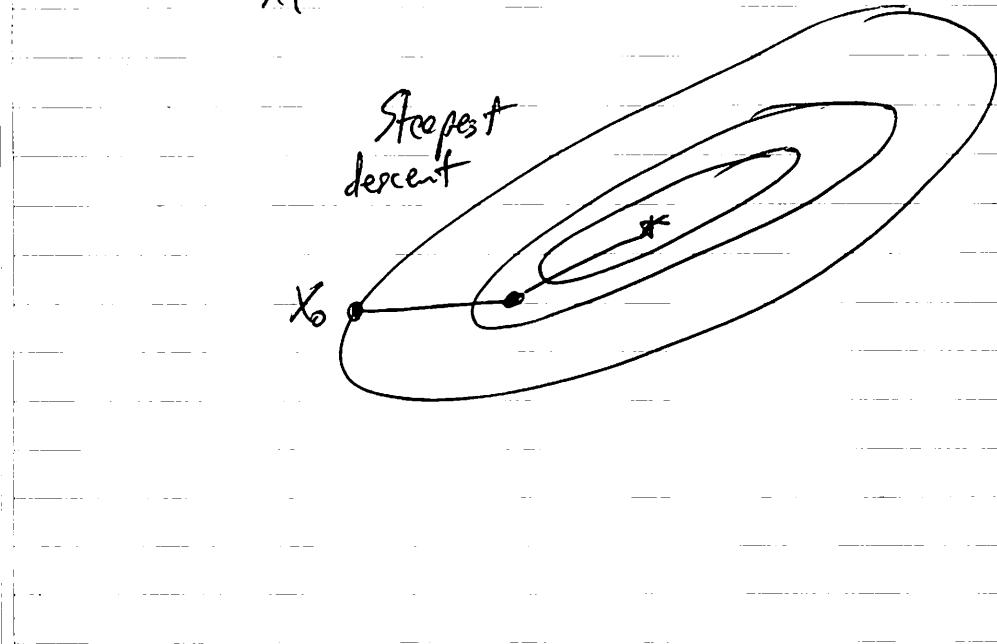
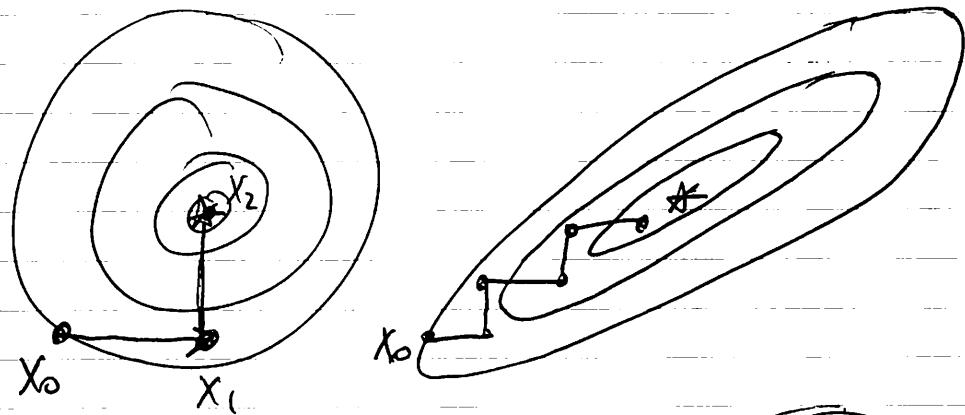
- ⇒ Method of alternating variables
- ⇒ cyclic coordinate descent
- ⇒ "deterministic Gibbs sampling"
- ⇒ backfitting

Coord. descent

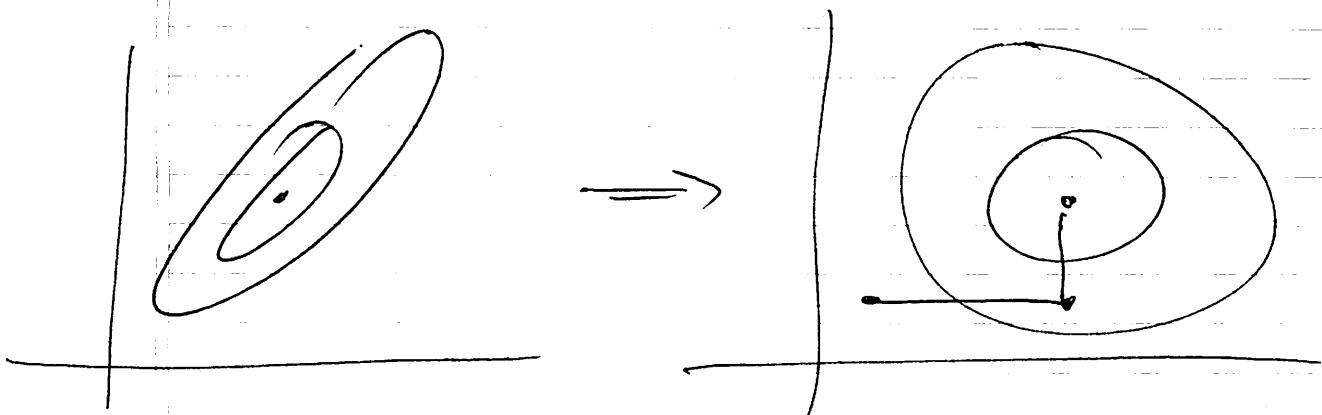
Coordinate descent can be very slow (slower than steepest descent).

But:

- ① Does not require calc. of f' or f''
- ② Often easier to do many 1-D mins than one k-D min.
- ③ Convergence can be good if coordinates are loosely coupled



Conjugate Gradient



Evaluate $f_0 = f(x_0)$, $f'_0 = f'(x_0)$

Let $p_0 = -f'(x_0)$

① Find $\min_{\alpha > 0} f(x_n + \alpha p_n) \Rightarrow \alpha_n$

Set $x_{n+1} = x_n + \alpha_n p_n$

② Eval $f'(x_{n+1}) = f'_{n+1}$

③ Let $\beta_{n+1} = \frac{f'^T f'_{n+1}}{f'^T f'_n}$

[Fletcher
Reeves]

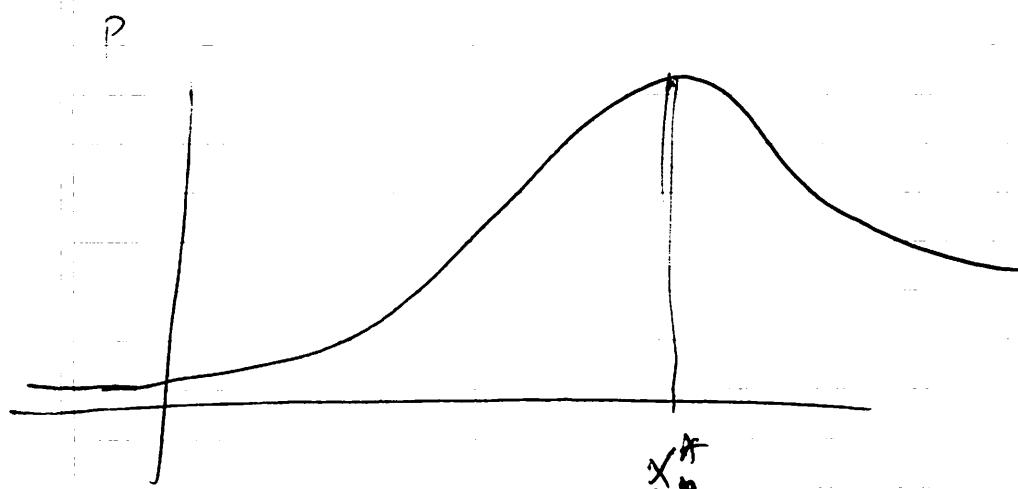
④ $p_{n+1} = -f'_{n+1} + \beta_{n+1} p_n$

Polak-Ribiere :

$$\beta_{k+1} = \frac{f'_{n+1}^T (f'_{n+1} - f'_n)}{f'_n^T f_n}$$

$$p_0 = -f'(x_0) = f'_0$$

$$p_1 = -f'_1 + \frac{f'_1^T f'_1 (-f'_0)}{f'_0^T f_0}$$



$$f(x) = \alpha x^2$$

$$f'(x) = 2x \quad f' = \begin{pmatrix} 2x + y \\ 2y + x \end{pmatrix}$$

$$p_0 = -2x_0 \Rightarrow x_1 = x_0 + \alpha (-2x_0)$$

$$p_1 = -2x_0 + \frac{4x_1^2}{4x_0^2} (-2x_0)$$

$$= -2x_1 + \frac{x_1^2}{x_0^2} (-2x_0)$$

Variations of Newton's Method

Again, solve $\ell'(\theta) = 0$

Newton's method: $\theta_{n+1} = \theta_n - \ell''(\theta_n)^{-1} \ell'(\theta_n)$

In general: $\theta_{n+1} = \theta_n - B_n^{-1} \ell'(\theta_n)$

(Fisher) scoring \Rightarrow replace ℓ'' , the observed information
with the expected information matrix.
(sometimes easier to compute)

Used to fit GLMs where it is equiv. to IRLS (with canonical link fn)

Quasi-Newton: Replace ℓ'' with a "seant-like" approximation

$$\textcircled{*} \quad \ell'(\theta_n) - \ell'(\theta_{n-1}) = B_n (\theta_n - \theta_{n-1})$$

Solve $\textcircled{*}$ for B_n (not unique, many ways).

Popular method due to Broyden, Fletcher, Goldfarb, and Shanno (BFGS). Also DFP (Davidon, Fletcher, Powell)

\Rightarrow In 1-D case, there is a unique solution

~~A Seant~~

"~~Second~~ Derivative"

unlike $-\ell''(\hat{\theta})$,
~~B_n~~ is not a
valid estimate
of $\text{Var}(\hat{\theta})$!!

$$\underbrace{l'(\theta_n) - l'(\theta_{n-1})}_{Y_n} = \underbrace{B_n (\theta_n - \theta_{n-1})}_{S_n}$$

$$B_n S_n = Y_n \quad \text{"secant equation"}$$

~~DPP~~ \rightarrow infinite solutions

$$\del{B_n = P^{-1} R^{-1} Y_n}$$

Add'l constraint: Find B closest to previous one,
and symmetric.

$$B_n = \arg \min_B \| B - B_{n-1} \| \text{ subj. to.}$$

$$B = B^T \text{ and } B S_n = Y_n$$

\Rightarrow solution is DFP method

$$\text{let } H_n = B_n^{-1}$$

$$\text{Solve } \min_H \| H - H_{n-1} \|$$

$$\text{subj. to } H = H^T \text{ and } H Y_n = S_n$$

\Rightarrow solution is BFGS method

Step-length selection

$$\begin{aligned} f'(x) &= f(x) \\ f'(x) &= f(x) \\ \phi'(0) &= f(x_0 + \alpha p_n) \\ &\approx f(x_0) + \alpha f'(x_0) p_n \end{aligned}$$

Given a step direction p_n , how far to go?

Let $\phi(\alpha) = f(x_0 + \alpha p_n)$.

Find $\min_{\alpha > 0} \phi(\alpha)$

\Rightarrow Too hard!

Roughly speaking:

① Choose initial α_0 . If

② $\phi(\alpha_0) \leq \phi(0) + c_1 \alpha_0 \phi'(0)$

then stop. ($c_1 \approx 10^{-4}$)

sufficient decrease condition

③ Otherwise, make quadratic approximation to $\phi(\alpha)$ called ϕ_q and let α_1 be the minimizer of ϕ_q .

If α_1 satisfies ②, stop.

④ Otherwise, make cubic approximation to ϕ , call ϕ_c , and let α_2 minimize ϕ_c .

If α_2 satisfies ②, stop

else repeat ③.

Cubic functions are good for approximating functions with much curvature.

`optimize()` in R uses polynomial (cubic) approximation.

~~Simulated annealing~~ → more later

Simplex Method

It's clear that in minimizing f , there is a trade off b/w knowledge of f and rate of convergence

Only Knowledge f only	partial Method	CCD (simulated annealing, simplex)	Sublinear (?) (linear)
f' only	steepest descent	linear	
partial f''	Quasi-Newton Fisher Scoring	super linear	
FULL f''	Newton	quadratic	fast



~~The EM Algorithm~~

$$\eta = \log \mu$$

Fisher Scoring - Poisson regression

$$Y_i \sim \text{Poisson}(\mu_i)$$

$$Y_i \sim \text{Poisson}(\exp(X^T \beta))$$

$$\frac{d\eta}{d\mu}$$

$$\frac{dY}{d\mu} = \frac{1}{\mu} \quad \eta_i = \log \mu_i = X_i^T \beta$$

$$\eta = X^T \beta$$

$$V(\mu) = \mu$$

$$\mu = \exp(X^T \beta)$$

① Start with $\hat{\mu}_0$

② Set $\hat{\eta} = \hat{\beta} + (Y - \hat{\mu}) / \hat{\mu}$

(adjusted response, working response)

$$\hat{\eta} = \log \hat{\mu}$$

Fisher Scoring for GLMs

$$Y \sim \text{Gaussian}(\mu), P(Y|\mu) \quad E[Y] = \mu \\ g(\mu) = X\beta \quad V(\mu) = V(\mu)$$

$$g(y) = g(\mu) + (y - \mu)g'(\mu)$$

\downarrow
z (working response)

$$\eta = g(\mu)$$

① Start with $\hat{\mu}_i$

② Set $z_i = \eta_i + (y_i - \hat{\mu}_i) \left(\frac{\partial \eta_i}{\partial \mu_i} \right)^*$

③ Weighted regression of z on X , i.e. solve

$$X^T W X \beta_n = X^T W z$$

$$\beta_n = (X^T W X)^{-1} X^T W z$$

where

$$W = \left[\left(\frac{\partial \eta_i}{\partial \mu_i} \right)^2 V(\mu_i) \right]^{-1}$$

e.g.: $y_i \sim \text{Poisson}(\mu_i)$

$$g(\mu_i) = \log \mu_i = \eta_i$$

$$\frac{\partial \eta_i}{\partial \mu_i} = \frac{1}{\mu_i} \quad V(\mu_i) = \mu_i$$

① set $\hat{\mu}_i$:

② $z_i = \log \mu_i + (y_i - \mu_i)/\mu_i$

③ Regress z on X where

$$w = \left[\frac{1}{\hat{\mu}_i^2} \hat{\mu}_i \right]^{-1} = \hat{\mu}_i$$

Poisson regression via Newton

$$L(\beta) \propto e^{-\mu} \mu^Y$$

$$l(\beta) = \gamma \log \mu - \mu$$

$$= Y \times p - \exp(X_p)$$

$$l'(\beta) = x'y - x' \exp(x\beta)$$

$$\ell''(\beta) = -x' \underbrace{\exp(x\beta)}_M x$$

$$\left(Y^T X \beta - \exp(X\beta) \right)$$

$$X^T y = X^T \mu$$

$$-\mathbf{x}^T \mathbf{w} \mathbf{x}$$

$$z = \eta + (\gamma - \mu) / \mu$$

$$g_{\mu z} \approx \mu \eta + \gamma^{-\mu}$$

$$\gamma - \mu^3 = \mu^2 - \mu^4$$

$$\eta = g(\mu) = \log \mu$$

$$\eta = e^{-\frac{1-\lambda}{\lambda}}$$

X β

$$\mu = e^{\lambda p}$$

1. =

$$\beta_{n+1} = \beta_n + (-x^T w x)^{-1} (-x^T (y - \mu))$$

$$= \beta_n + (X^T W X)^{-1} X^T (\mu z - \mu y)$$

$$= \beta_n + (X^T W_n X)^{-1} X^T W_n (z_n - \eta_n)$$

$$= (X^T W X)^{-1} X^T W Z + \beta_0 - (X$$

$$= (X^T W X)^{-1} X^T W z + \underbrace{(\beta_n - (X^T W X)^{-1} X^T W y)}_{V}$$

$$= (x^T W x)^{-1} x^T W z$$

二〇

$$= \log M$$

$$= x\beta$$

The EM algorithm

~~EM stands for Expectation Maximization.~~

~~Originally outlined by DLR (1977) but may go much farther back. DLR worked many different ideas and gave examples of the broad applicability of EM.~~

~~EM is not an "algorithm". It is an algorithm for creating other algorithms.~~

Generalized Additive Models (Hastie + Tibshirani)

Usual linear model has:

$$\text{Fits } Y_i = \alpha + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_p X_{pi} + \varepsilon_i$$

$$Y_i = \alpha + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_p X_{pi} + \varepsilon_i$$

GAMs say:

$$Y_i = \alpha + s_1(X_{1i}) + s_2(X_{2i}) + \dots + s_p(X_{pi}) + \varepsilon_i$$

$$\Rightarrow Y = \alpha + \sum_{j=1}^p s_j(X_j) + \varepsilon \quad (\#)$$

where $s_j()$ are smooth.

$s_j()$ can be any kind of "smoother", even a mixture of different kinds

ex. Smoothers: parametric splines, smoothing splines, penalized splines, loess, running lines, running median

GAM algorithm: $\rightarrow Y = (Y_1, \dots, Y_n)$

Given model $Y = \alpha + \sum_{j=1}^P S_j(X_j) + \varepsilon$

(1) Initialize $\alpha = \frac{1}{n} \sum_{i=1}^n Y_i$, $s_j = s_j^0 = 0$

(2) For $j = 1, 2, \dots, P$

Let $r_j = Y - \alpha - \sum_{k \neq j} S_k(X_k)$

where ~~$S_k = (S_k(x_{k1}), S_k(x_{k2}), \dots, S_k(x_{kn}))$~~

where $S_k = (S_k(x_{k1}), S_k(x_{k2}), \dots, S_k(x_{kn}))$

i.e. $S_k()$ evaluated at data points for X_k

$\hat{s}_j = \text{Smooth}(r_j | X_j)$

$\hat{s}_j = s_j - \sum_{i=1}^n s_j(x_{ij})$

(3) Evaluate $\Delta = \sum_{j=1}^P \|s_j - \hat{s}_j\|$

or $\Delta = \frac{\sum_{j=1}^P \|s_j - \hat{s}_j\|}{\sum \|s_j\|}$

If $\Delta < \varepsilon$, stop.

Else set $s_j^* = \hat{s}_j$ for $j = 1 \rightarrow P$.

~~Go to (2)~~

Local Scoring

Set adjusted response

$$z_i = \eta_i + (y_i - \mu_i) \frac{\partial \eta_i}{\partial \mu_i}$$

Fit an additive model for

$z \sim x_1 \rightarrow x_p$ with observation weights

$$W_i = \left[\frac{\partial \eta_i}{\partial \mu_i} \right]^{-2} V(\mu_i)^{-1}$$

GAM algorithm = "backfitting"

~ Alternating conditional expectation

~ Cyclic coordinate descent

⇒ linear convergence algorithm

⇒ sidesteps curse of dimensionality by additivity constraint (big!)

Boosting

Outcomes are $y_i \in \{-1, 1\}$ with covariates x_i :

Given $g(x) \in \mathbb{R}$, then let

$$f(x) = \text{sign}(g(x)) \in \{-1, 1\}$$

e.g. logistic regression

$$g(x) = \logit \frac{\mu(x)}{1 - \mu(x)} = x' \beta \in (-\infty, \infty)$$

$$\Rightarrow f(x) = \text{sign}(x' \beta)$$

Margin : $y_i f(x_i)$

$y_i f(x_i) > 0$ good!

$y_i f(x_i) < 0$ bad!

Loss function: Classifiers minimize some loss function

Binomial deviance:

$$L(y, f(x)) = \log(1 + e^{-2yf(x)})$$

Squared error:

$$L(y, f(x)) = (y - f(x))^2$$

Exponential:

$$L(y, f(x)) = e^{-yf(x)}$$

Ada Boost

- ① Set $w_i = \frac{1}{n}$ for $i=1 \rightarrow n$
- ② For $m = 1 \rightarrow M$
 - a) Fit classifier $f_m(x)$ to training data weighted by w_i
 - b) Compute
$$err_m = \frac{\sum_{i=1}^n w_i \mathbb{1}\{y_i \neq f(x_i)\}}{\sum_{i=1}^n w_i}$$
 - c) Let $c_m = \log \frac{1-err_m}{err_m}$
 - d) Update $w_i = w_i \exp(c_m \mathbb{1}\{y_i \neq f(x_i)\})$ for $i=1 \rightarrow n$
- ③ Final classifier, $\exists F(x) = \text{Sign} \left[\sum_{m=1}^M f_m(x) \right]$

Boosting

Friedman, Hastie, Tibshirani
2000 *Annals*

① Ada Boost \rightarrow why this?

$$\mathbb{E}[e^{-yF(x)}] \rightarrow \text{minimized by } \frac{1}{2} \log \frac{\mathbb{P}(y=1|x)}{\mathbb{P}(y=-1|x)}$$

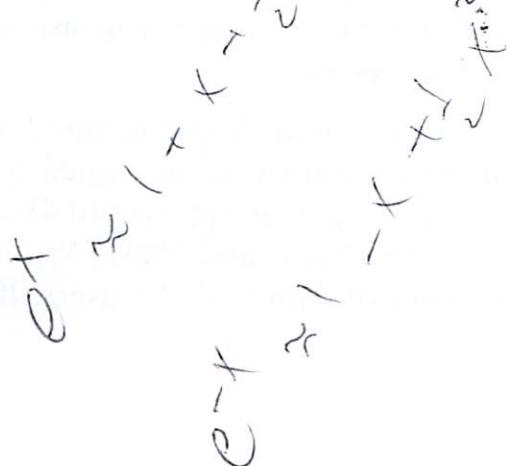
PF:

$$\frac{\partial}{\partial F} \mathbb{E}[e^{-yF(x)}|x] = -e^{-F(x)} \mathbb{P}(y=1|x) + e^{F(x)} \mathbb{P}(y=-1|x) = 0$$

$$-e^{-F(x)} \mathbb{P}(y=1|x) = -e^{F(x)} \mathbb{P}(y=1|x)$$

$$\frac{\mathbb{P}(y=1|x)}{\mathbb{P}(y=-1|x)} = \frac{-e^{F(x)}}{-e^{-F(x)}} = e^{2F(x)}$$

$$\frac{1}{2} \log \frac{\mathbb{P}(y=1|x)}{\mathbb{P}(y=-1|x)} = F(x)$$



② AdaBoost \Rightarrow Additive logistic regression model
 via Newton-like steps
 for minimizing $\mathbb{E}[e^{-YF(x)}]$

PF: Let $J(F) = \mathbb{E}[e^{-YF(x)}]$. Given x and c , improve $F(x)$ by taking a step $F(x) + c f(x)$.

\Rightarrow Newton's Method was $X_{n+1} = X_n + \underbrace{f''(x_n)^{-1} f'(x_n)}_{\substack{\downarrow \text{function} \\ \downarrow \text{scalar}}} \quad \left. \begin{array}{l} \\ = X_n + (1) f(x_n) \end{array} \right\}$

$$J(F + cf) = \mathbb{E}[e^{-Y(F+cf)}] = \mathbb{E}[e^{-YF} e^{-Ycf}]$$

$$\approx \mathbb{E}[e^{-YF}(1 - Ycf + \frac{1}{2}Y^2c^2f^2)]$$

$$\begin{aligned} Y^2 &= 1 \\ f^2 &= 1 \end{aligned} \quad = \mathbb{E}[e^{-YF}(1 - Ycf + \frac{1}{2}c^2)]$$

$$\text{We want } f = \underset{f}{\operatorname{arg\,min}} \mathbb{E}[e^{-YF}(1 - Ycf + \frac{1}{2}c^2)]$$

$$\text{for } \underset{x}{\operatorname{arg\,min}} \underset{f}{\operatorname{arg\,min}} \mathbb{E}_w[1 - Ycf + \frac{1}{2}c^2]$$

$$H \stackrel{1}{=} (Y^2 - 2Yf + f^2) - 1$$

$$\text{P}(\cdot) = -f_Y + 2f$$

For $C > 0$,

$$\text{min } \mathbb{E}_w[-Yf + \frac{1}{2}C^2]$$

\Leftrightarrow

$$\max \mathbb{E}_w[Yf] = \frac{\mathbb{E}[e^{-YF} Yf | x]}{\mathbb{E}[e^{-YF} | x]}$$

$$\cancel{\mathbb{E}_w[Yf]} = \mathbb{E}_w[Y | x] \quad -\mathbb{E}[Yf] = \mathbb{E}[(Y-f)^2] \frac{1}{2} - 1$$

$$\cancel{e^{-YF} P(Y=1)} = e^{-F} P(Y=1)$$

$$\cancel{e^{-F} P(Y=-1)} = e^{F} P(Y=-1) = \cancel{\frac{1}{2}} (1 - \cancel{e^{-F}})^2$$

$$= f(x) \mathbb{E}_w[Y | x]$$

$$= f(x) \left[1 P_w(Y=1 | x) - 1 P_w(Y=-1 | x) \right]$$

$$= f(x) \left[P_w(Y=1 | x) - P_w(Y=-1 | x) \right]$$

$$f(x) = \begin{cases} 1 & \text{if } P_w(Y=1 | x) \geq P_w(Y=-1 | x) \\ -1 & \text{otherwise} \end{cases}$$

β Minimize

To get c , we minimize J for fixed $f(x)$

$$J(F+cf) = \mathbb{E}[e^{-YF - cf}]$$

$$= \mathbb{E}_w[e^{-Ycf}]$$

$$= e^{-cf} \mathbb{E}_w[1\{Y=1\}] + e^{cf} \mathbb{E}_w[1\{Y=-1\}]$$

$$\frac{\partial}{\partial c} J(F+cf) = -f e^{-cf} \mathbb{E}_w[1\{Y=1\}] + f e^{cf} \mathbb{E}_w[1\{Y=-1\}] = 0$$

$$e^{cf} \mathbb{E}_w[1\{Y=-1\}] = e^{-cf} \mathbb{E}_w[1\{Y=1\}]$$

$$e^{2cf} = \frac{\mathbb{E}_w[1\{Y=1\}]}{\mathbb{E}_w[1\{Y=-1\}]}$$

$$cf = \frac{1}{2} \log \frac{\mathbb{E}_w[1\{Y=1\}]}{\mathbb{E}_w[1\{Y=-1\}]}$$

$$\Rightarrow c = \frac{1}{2} \log \frac{\mathbb{E}_w[1\{Y=1\}]}{\mathbb{E}_w[1\{Y=-1\}]} = \frac{1}{2} \log \frac{1 - er}{er}$$

$$er = \mathbb{E}_w[1\{Y \neq f(x)\}]$$

The new $F(x)$ is the

$$F_{\alpha}(x) = F(x) + c f(x)$$

$$= F(x) + \frac{1}{2} \log \frac{1 - \epsilon_{\text{err}}}{\epsilon_{\text{err}}} f(x)$$

New weights are

$$e^{-yF} \Rightarrow e^{-yF - cf(x)y}$$

because $-yf(x) = 2 \mathbb{1}\{y \neq f(x)\} - 1$, new weights are

$$e^{-yF} \left(\exp \left(\frac{1}{2} \log \frac{1 - \epsilon_{\text{err}}}{\epsilon_{\text{err}}} (2 \mathbb{1}\{y \neq f(x)\} - 1) \right) \right)$$

$$= e^{-yF} \left[\# \exp \left(\log \frac{1 - \epsilon_{\text{err}}}{\epsilon_{\text{err}}} \mathbb{1}\{y \neq f(x)\} \right) \right] (k)$$

\Rightarrow Same as Discrete Adaboost

Biostat 778: Homework 1

November 16, 2016

1 Improving consistent estimators

Let $\tilde{\theta}_n$ be an estimator of θ such that $\sqrt{n}(\tilde{\theta}_n - \theta) \rightarrow \mathcal{N}(0, \Sigma(\theta))$ and $\Sigma(\theta) < \infty$. Let $\hat{\theta}_n$ be the maximum likelihood estimator of θ . Let $\tilde{\theta}_n^{(1)}$ be single iteration of Newton's method applied to $\tilde{\theta}$, i.e.

$$\tilde{\theta}_n^{(1)} = \tilde{\theta}_n - \ell''(\tilde{\theta}_n)^{-1}\ell'(\tilde{\theta}_n)$$

Show that $\tilde{\theta}_n^{(1)}$ is asymptotically equivalent to the MLE.

2 Logistic Regression with Penalization

Write a function that can fit a logistic regression model while allowing for an L^2 penalty on the parameters. That is, if $\ell(\beta)$ is the log-likelihood for the parameter vector β , then you want to maximize the penalized log-likelihood

$$\ell^*(\beta, \lambda) = \ell(\beta) - \lambda\beta'\beta$$

Write your own implementation of Newton's method to optimize the penalized likelihood. The output of the function should be the maximum (penalized) likelihood estimates of β and the asymptotic standard errors for each of the elements of $\hat{\beta}$.

Summary of Minimization

$$\min_x f(x) \quad \text{for} \quad f: \mathbb{R}^k \rightarrow \mathbb{R}, x \in \mathbb{R}^k$$

Line Search: ~~Step~~

① Steepest descent

$$x_{n+1} = x_n + \alpha [-f'(x_n)]$$

↓
direction of steepest descent
scalar step length

⇒ linear convergence

② Newton

$$x_{n+1} = x_n + \underbrace{\left[f''(x_n) \right]^{-1}}_{\text{Newton direction}} \underbrace{[-f'(x_n)]}_{\text{Step length } = 1}$$

Step length = 1

⇒ quadratic convergence

→ in stat, estimate of asymptotic covariance

③ Quasi-Newton

$$x_{n+1} = x_n + \alpha B_n [-f'(x_n)]$$

$$\textcircled{a} \quad B_n = \arg \min_B \|B - B_{n-1}\|$$

$$B = B^T \quad \text{and} \quad f'(x_n) - f'(x_{n-1}) = B_n (x_n - x_{n-1})$$

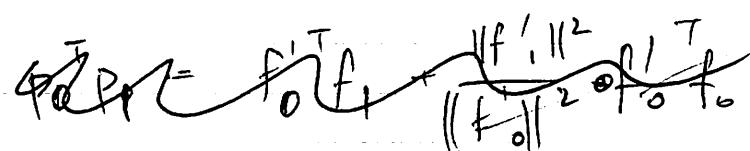
④ Conjugate Gradient ("Modified steepest descent")

$$x_{n+1} = x_n + \alpha p_n$$



$$p_{n+1} = -f'_n(x_{n+1}) + \frac{\|f'(x_{n+1})\|^2}{\|f'(x_n)\|^2} p_n$$

$$p_0 = -f'(x_n)$$



Coordinate Descent:

$$x_{n+1}^{(k)} = \arg \min_{x_n^{(k)}} f(x_n^{(1)}, x_n^{(2)}, \dots, x_n^{(k)}, \dots, x_n^{(P)})$$

\Rightarrow has "descent property", non-increasing

\Rightarrow thus

~~Method of steepest descent / ascent uses~~
 ~~$On+1 = On - \alpha \cdot l'(On)$~~
~~Where α is chosen so that $On+1$ has a larger likelihood value than On .~~

~~Knowledge of l''~~ Summary:



① Newton: Fastest (quadratic convergence), requires calculating l'' , gives asymptotic $\text{Var}(\hat{\theta})$, unstable if starting value too far.

~~

② Scoring: Super linear convergence, but equiv. to Newton, if l'' does not depend on θ (the many common cases), requires calculating $E[l'']$, can be unstable, but often quite stable in typical statistical apps.

~~

③ Quasi-Newton: Super linear convergence, does not require l'' or $E[l'']$, more stable than Newton, No estimate of $\text{Var}(\hat{\theta})$

④ ~~Steepest descent~~: Linear convergence, very slow

✗

④ Steepest descent: linear convergence, stable

Secant method

If f' is difficult to compute (almost always) or we are lazy (always) the secant method provides an approximation.

Newton step / Secant method

$$\begin{aligned} x_n - x_{n-1} &= \\ x_n - x_\infty + x_\infty - x_{n-1} &= \\ \varepsilon_n &\rightarrow \varepsilon_{n-1} \end{aligned}$$

$$x_{n+1} = \cancel{x_n - f(x_n)}$$

$$x_{n+1} = x_n - \frac{f(x_n)(x_n - x_{n-1})}{f(x_n) - f(x_{n-1})}$$

$$\Rightarrow f'(x_n) \approx \frac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}}$$

~~PRO~~ Note: Need 2 starting values.

$$\varepsilon_{n+1} = \varepsilon_n -$$

Pro: Easy as pie

$$\frac{f(x_\infty + \varepsilon_n)(x_n - x_{n-1})}{f(x_\infty + \varepsilon_n) - f(x_\infty + \varepsilon_{n-1})}$$

Con: Convergence only super linear.



Statistics application

$$\left(\frac{\varepsilon_n - \frac{\varepsilon_n^2}{2} f''}{\frac{\varepsilon_n f' + \frac{\varepsilon_n^2}{2} f''}{\varepsilon_n f' - \frac{\varepsilon_n^2}{2} f''} (\varepsilon_n - \varepsilon_{n-1}) f' - \left(\frac{\varepsilon_n^2 - \varepsilon_{n-1}^2}{2} \right) f''} \right)$$

$$\begin{aligned} &= \varepsilon_n - \\ &\quad \frac{\left(\varepsilon_n f' + \frac{\varepsilon_n^2}{2} f'' \right) (\varepsilon_n - \varepsilon_{n-1})}{\left(\varepsilon_n - \varepsilon_{n-1} \right) f' - \left(\frac{\varepsilon_n^2 - \varepsilon_{n-1}^2}{2} \right) f''} \end{aligned}$$

$\ell(\theta): \mathbb{R}^K \rightarrow \mathbb{R}$, log-likel.hood, $\theta = (\theta_1, \dots, \theta_K)$

$\ell'(\theta): \mathbb{R}^K \rightarrow \mathbb{R}^K$, ~~gradient~~ gradient

$\ell''(\theta): \mathbb{R}^K \rightarrow \mathbb{R}^{K \times K}$, hessian