① Solutions to non linear equations

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

- Bisection

- Functional iteration

- Newton's method

Stat Model → Technique + principle (← Data, → Statistic) → Algorithm → Program

All Statistics is:

① Probability
② Linear Algebra
③ Optimization ⟸

We look at ③.

Generally, we want to maximize or minimize something.

⟹ max! likelihood
⟹ min! sum of squares

⟹ ~~max~~ max $f$ = min $-f$
So don't worry about it.

This course is about

① maximizing a function
② integrating a function

Deterministic Algorithms

General Idea:

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 $g$.

④ Iterate ②, ③

$\Rightarrow$ Instead of direct max, "transfer" to simpler function and iterate

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

$$f(x_{n-1}) - f(x_{\infty}) = f'(c)(x_{n-1} - x_{\infty})$$

$$f(x_{n-1}) = f'(c)(x_{n-1} - x_{\infty})$$

~~Course Outline~~

Course outline

Maximizing

(★) Solving non linear equations (root finding)

$$f(x) = 0 \qquad \text{for } f: \mathbb{R} \to \mathbb{R}$$

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

$$K = 2, 3, \ldots, N \quad (\text{~~~~~})$$

(★) General optimization routines

$$\text{Given } f: \mathbb{R}^K \to \mathbb{R}, \quad \underset{x}{\text{Max}} \, f(x)$$

$$\text{or } \underset{x}{\text{min}} \, f(x)$$

Line search methods         | Taylor's theorem
  ↳ Newton
  ↳ Quasi — Newton

① Pick a direction
② Go a certain distance $M$ that direction

related: Simulate annealing

A random optimizer — still general purpose

(★) Statistics !

EM algorithm for maximum likelihood

Minorization / majorization

Monte Carlo EM

densities: $p(x) = C f(x)$

↑ integral

Course outline

Integration "How to compute an integral"

Miscellaneous

* 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

# Solving non-linear equations

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

## Bisection method.

$$\text{If} \quad \text{sgn}\left(f(a)\right) \neq \text{sgn}\left(f(b)\right)$$

$\Rightarrow$ Intermediate value thm

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

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

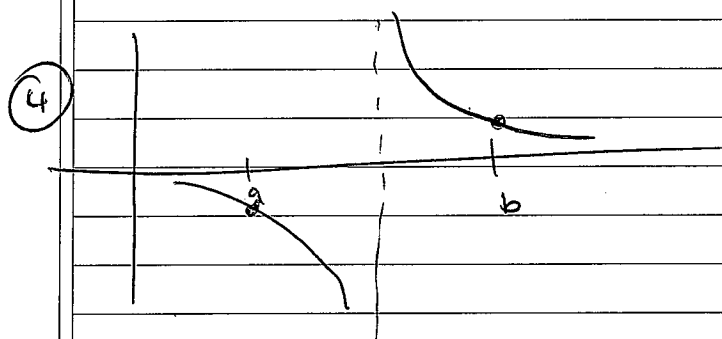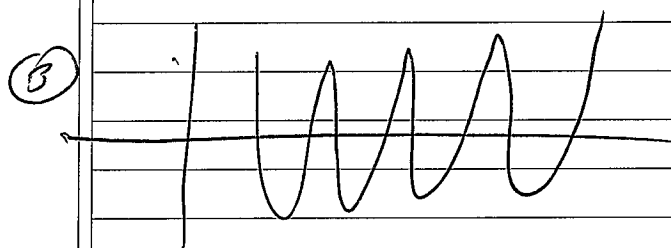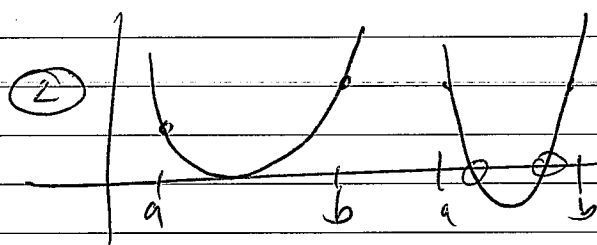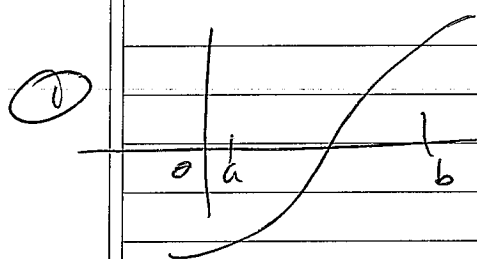① Let $c = \dfrac{a+b}{2}$

② If $f(c) = 0$, stop

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

④ goto ①

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

① 

② 

③ 

④ 

Converge when $|b - a| < \varepsilon$ or

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

Ex.

$l(\theta) = $ likelihood

$l'(\theta) = 0 \implies \hat{\theta}$ 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 $K$,

$$S_k \overset{\Delta}{=} \{\theta : f(\theta \mid y) \geq K\}$$



$(a \qquad b]$

## Ex. likelihood intervals

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

find $\theta$. LI $= \{\theta : f(\theta) \geq 1/8\}$

Solve $f(\theta) - 1/8 = 0$

Ex. Bayesian credible interval

Let $S_k = \{\theta : f(\theta \mid y) \geq K\}$

Bayesian credible interval of level $\alpha$ $\approx$ finds $K$ s.t.

$$\mathbb{P}(\theta \in S_K \mid y) = \alpha$$

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

$\mu([a,b]_k) - \alpha = 0$ \qquad Solve for $K$

$[a, b]_k$ \qquad\qquad Solve for $a, b$.

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

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

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

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

Bisection algorithm: is interval length $\propto \frac{1}{2^n}$

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

$\vdots$

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

# Rates of convergence

Suppose $X_n \to X_\infty^*$ in $\mathbb{R}^k$. ~~Then we have~~

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

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

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

$$X_\infty = 1$$

~~Q-Quadratic~~
Q-superlinear if

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

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

Q-Quadratic if

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

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

- test vectors
- homework/reading
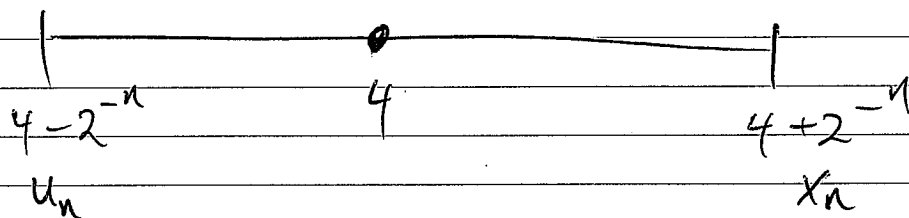- office hours
- no TA

① a) $X_n = 4 + 2^{-n}$

$$\frac{||X_{n+1}||}{||X_n||} = \frac{4 + 2^{-(n+1)}}{4 + 2^{-n}} = \frac{4}{4 + 2^{-n}} + \frac{2^{-(n+1)}}{4 + 2^{-n}}$$

$$< 1$$

② b) $X_n = 10000 + 2^{-n}$

$$\frac{X_{n+1}}{X_n} = \frac{1000}{1000 + 2^{-n}} + \frac{2^{-(n+1)}}{1000 + 2^{-n}}$$

$$\approx 1 + 2 - 0$$

Want to know ratio of errors, not ratio of sequence elements



$$4 - 2^{-n} \qquad 4 \qquad 4 + 2^{-n}$$
$$u_n \qquad\qquad\qquad X_n$$

$u_n$ and $X_n \to 4$ at same rate, but

$$\frac{X_{n+1}}{X_n} < 1 \qquad \text{and} \qquad \frac{|u_{n+1}|}{|u_n|} > 1$$

$$\frac{4 - 2^{-(n+1)}}{4 - 2^{-n}}$$

Ex. A bisection algorithm:

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

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

$$= 2^{-n-1+n} = \frac{1}{2} < r \in (0,1)$$

Bisection achieves linear convergence

~~Quest Newt~~

Newton's Method    —    quadratic
Quasi - Newton     =    superlinear
steepest descent   —    linear
~~Bisection~~       o

Functional Iteration

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

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 $X_n \to X_\infty$ where $f(X_\infty) = X_\infty$ (i.e. a fixed point).

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

When does functional iteration work?

### Newton's Method

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

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

$z$ is b/w $X_n$ and $X_\infty$.

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

Since $X_\infty$ and $z$ unknown, do

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

$$\Downarrow$$

Newton update

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

## Shrinking Lemma

Let $M$ be a closed subset of a c.n.v.s. Let $f : M \to M$ 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 $x_0$

Proof:

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

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

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

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

Left margin calculations:

$$S = 1 + \frac{1}{r} + \frac{1}{r^2} + \cdots$$
$$rS = r - 1 + \frac{1}{r} - \frac{1}{r^2} + \cdots$$
$$rS = r + S$$
$$r(S-1) = S$$
$$r = \frac{S}{S-1}$$
$$S(r-1) = r$$
$$S = \frac{r}{r-1} = \frac{1}{1-\frac{1}{r}}$$

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

$$\sum_{k=0}^{\infty} \frac{1}{r^k} = 1 + \frac{1}{r} + \frac{1}{r^2} + \cdots$$
$$a = 1 + \frac{1}{r} + \frac{1}{r^2} + \cdots$$
$$a = r\left(\frac{1}{r} + \frac{1}{r^2} + \cdots\right)$$

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

$f^m$

---

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

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

$$\le \frac{1}{1-K} \left| f(x) - x \right|$$

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

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

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

$\Rightarrow \{ f^n(x) \}$ is a Cauchy sequence. Let $x_0$ be its limit.

Let $N$ be s.t. $\forall n \ge N$, $\left| f^n(x) - x_0 \right| < \varepsilon$.

Then

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

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

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

Let $x_1$ be another fixed point. Then

$$\left| x_1 - x_0 \right| = \left| f(x_1) - f(x_0) \right| \le K \left| x_1 - x_0 \right|.$$

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

Thm:

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

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

converges to $x_\infty$.

Proof:

Note that

$$g'(x) = 1 - \frac{f(x)f''(x) - f'(x)f'(x)}{[f'(x)]^2} \quad 1 - \frac{f'(x)f'(x) - f(x)f''(x)}{[f'(x)]^2}$$

$$= \frac{f(x)f''(x)}{[f'(x)]^2}$$

$$\Rightarrow g'(x_\infty) = 0$$

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

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

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

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

$$|g(a) - g(b)| \leq |g'(c)||a - b|$$

$$\leq k|a - b| \qquad (0 < k < 1)$$

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

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

$\text{Max } f = \text{classical / frequentist}$

$\int f \, d\mu = \text{Bayesian}$

Convergence rates for shrinking maps.

~~Suppose $g : \mathbb{R}^n \to \mathbb{R}$ and~~

Suppose $g$ satisfies

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

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

Also, assume $0 < |g'(x_\infty)| < 1$, where $x_\infty$ is the fixed point. Then $x_n \to x_\infty$ at a linear rate.

Pf: $\dfrac{|x_{n+1} - x_\infty|}{|x_n - x_\infty|} = \dfrac{|g(x_n) - g(x_\infty)|}{|x_n - x_\infty|}$

Taking limits

$$\lim_{n \to \infty} \frac{|g(x_n) - g(x_\infty)|}{|x_n - x_\infty|} = \underbrace{|g'(x_\infty)| > 0}_{\text{constant} \in (0,1)}$$

$\Rightarrow$ linear convergence

What about Newton's method?

~~For~~ Suppose $f \in C^2$, ~~and~~ and $\exists \, x_\infty$ s.t.

$f(x_\infty) = 0$.

By Taylor's theorem: for some small $\varepsilon$,

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

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

② $f'(x_\infty + \varepsilon) = f'(x_\infty) + \varepsilon f''(x_\infty) + O(\varepsilon)$

~~#~~ Newton's method generates the sequence

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

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

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

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

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

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

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

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

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

$$\sim \frac{f''(x_\infty)/2}{f'(x_\infty)} \qquad \varepsilon_n \downarrow 0$$

$$\Rightarrow \exists \text{ some } M < \infty \text{ s.t.}$$

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

$$\Rightarrow \text{Quadratic convergence.}$$

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

In practise, we ignore assumptions/conditions.
Use Newton's method as a "black box". Cavert Emptor.

Pro: Very fast in neighborhood of truth
Direct multivariate generalization

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

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

## Newton's method

$$\Theta_{n+1} = \Theta_n - \left[\ell''(\Theta_n)\right]^{-1} \ell'(\Theta_n)$$

$\qquad\qquad\quad\; \underbrace{\phantom{XX}}_{K \times 1} \quad \underbrace{\phantom{XXXX}}_{K \times K} \qquad \underbrace{\phantom{XX}}_{K \times 1}$

$\Rightarrow$ May be easier/better to solve

$$\underbrace{\left[\ell''(\Theta_n)\right]}_{A} \underbrace{\Theta_{n+1}}_{x} = \underbrace{\left[\ell''(\Theta_n)\right]\Theta_n - \ell'(\Theta_n)}_{b}$$

Than try to invert $\ell''(\Theta_n)$.

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

$\ell'(\hat{\Theta}) =$ score statistic

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

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

$$\sqrt{n}\left(\left[-\ell''(\hat{\Theta})^{-1/2}\right](\hat{\Theta} - \Theta_0)\right) \longrightarrow N(0, \mathcal{I})$$

for $n \longrightarrow \infty$.

Left margin notes:

$$\log P_i - (y_i - P_i = x_i^T \beta$$

$$\log y_i - x_i^T \beta$$

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

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

$$W = diag[P_i(1-P_i)]$$

Main body:

$$Y_i \sim Bernoulli(P_i), \quad i = 1, \to n$$

$$logit(P_i) = \log \frac{P_i}{1-P_i} = x_i^T \beta$$

$$P_i = \frac{e^{x_i^T \beta}}{1 + e^{x_i^T \beta}}$$

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

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

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

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

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

$$= \sum Y_i\left[x_i - x_i P_i\right] + (1-Y_i)(-x_i P_i)$$

$$(Y_i - 1)(x_i P_i)$$

$$= \sum Y_i\left[x_i - x_i P + x_i P_i\right] - x_i P_i$$

$$= \sum Y_i x_i - x_i P_i$$

$$= \sum x_i(Y_i - P_i) = X^T(Y-P)$$

$$\ell'(\beta) = X^T(y-p)$$

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

$$\hookrightarrow \quad \text{diag}\left[ p_i(1-p_i) \right]$$

$$\beta_{n+1} = \beta_n + \left[ -X^T W_n X \right]^{-1} \left[ X^T(y-p_n) \right]$$

Newton update

For exp. families w/ canonical link,

$\bullet\ \ell''(\theta) = \mathbb{E}\,\ell'(\theta)$ so Newton and Fisher scoring are same

# General Purpose Minimization

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

usually $\Longleftrightarrow f'(x) = 0$

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

$\hookrightarrow$ 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 is steepest-descent : ~~$f'(x_n)$~~

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

$\hookrightarrow$ orthogonal to contours of $f$

# Newton direction:

By Taylor's Theorem:

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

$$\Downarrow$$

$$m_n(p)$$

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

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

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

That's familiar!

Similarly, Quasi-Newton

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

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

$$\downarrow$$

$B_n$ satisfies a "secant condition"

# Coordinate descent

If $f$ is $K$-dimensional, we minimize along individual dimensions (coordinate) in a cyclic fashion.

$\Rightarrow$ Method of alternating variables
$\Rightarrow$ cyclic coordinate descent
$\Rightarrow$ "deterministic Gibbs sampling"
$\Rightarrow$ backfitting

# 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 $\ell''$, 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 fcn)

Quasi-Newton: Replace $\ell''$ with a "secant-like" approximation

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

Solve $\circledast$ 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

~~$\circledast$ $S_n = Y_n$     "Secant equation"~~

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

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

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

~~DFP~~ $\longrightarrow$ infinite solutions

~~$B_n = (I - \ldots)YY^T$~~

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

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

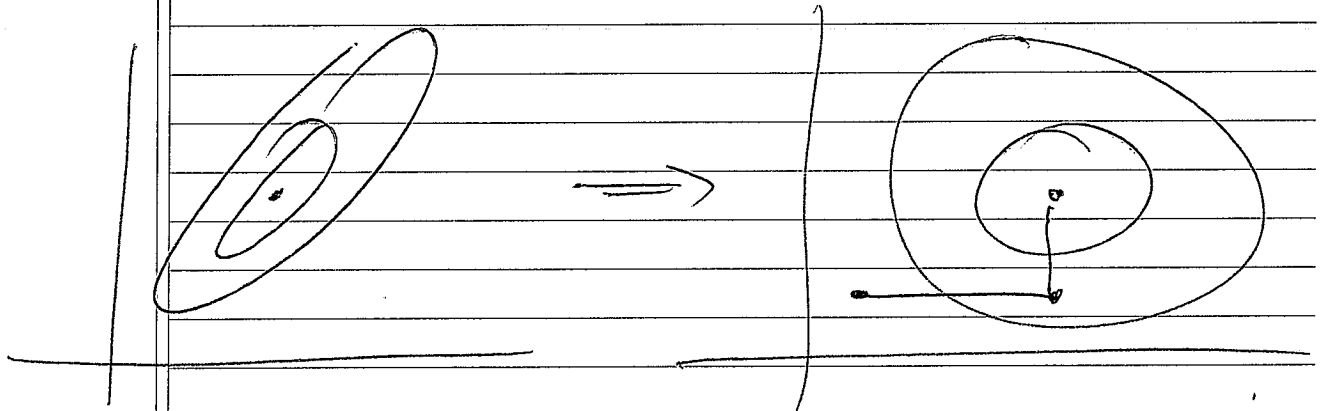$$B = B^T \quad \text{and} \quad B S_n = Y_n$$

$\Longrightarrow$ solution is DFP method

Let $H_n = B_n^{-1}$

Solve $\min_{H} \| H - H_{n-1} \|$

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

$\Longrightarrow$ solution is BFGS method

# Conjugate Gradient



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

Let $p_0 = -f'(x_0)$

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

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

②  Eval $f'(x_{n+1}) = f'_{n+1}$

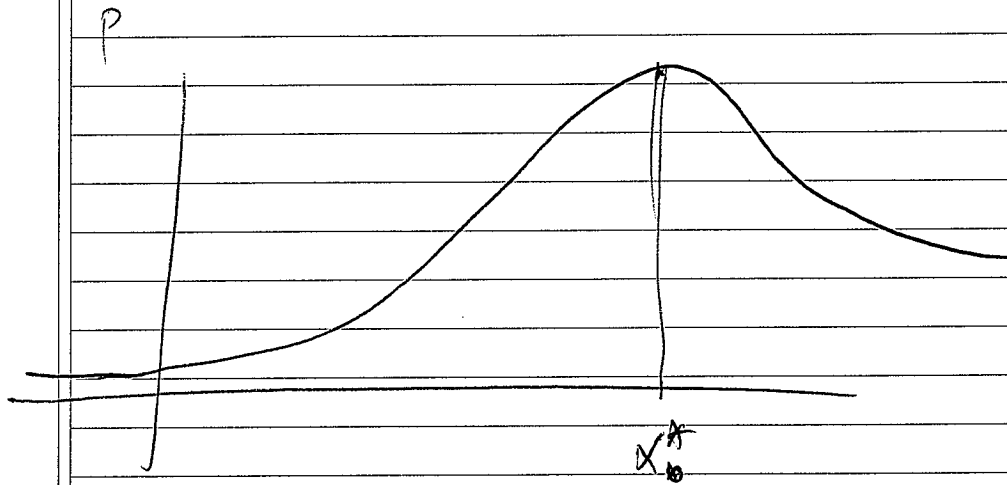③  let $\beta_{n+1} = \dfrac{f_{n+1}'^T f_{n+1}'}{f_n'^T f_n'}$    $\left[\begin{array}{c}\text{Fletcher}\\ \text{— Reeves}\end{array}\right.$

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

Polak - Ribiere:

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

$$p_0 = -f'(x_0) = -f_0'$$

$$p_1 = -f_1' + \frac{f_1'^T f_1'}{f_0'^T f_0}(-f_0')$$

$P$



$x_0^A$

$$f(x) = x^2 + xy - \cancel{xy}$$

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

$$p_0 = -2x_0 \qquad \Rightarrow \quad x_1 = x_0 + \alpha(-2x_0)$$
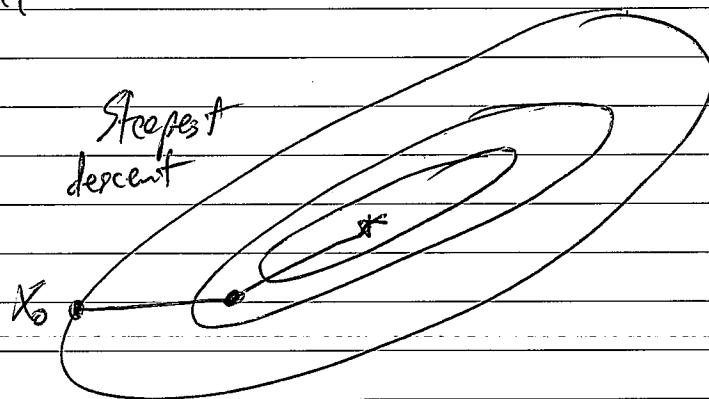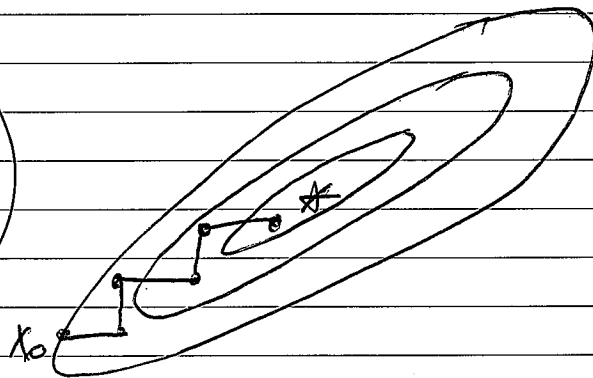
$$p_1 = -2x_0 + \frac{4x_1^2}{4x_0^2}(-2x_0)$$
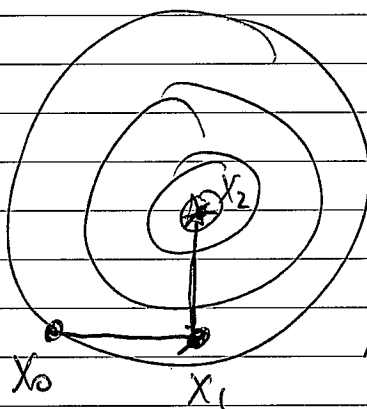
$$= -2x_1 + \frac{x_1^2}{x_0^2}(-2x_0)$$

## Coord. descent

<u>Coordinate descent</u> can be <u>very</u> 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



Steepest descent

# Step-length selection

Given a step direction $p_n$, how far to go?

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

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

$\Rightarrow$ Too hard!

Roughly speaking:

① Choose initial $\alpha_0$. If

Ⓐ $\quad \phi(\alpha_0) \le \phi(0) + c_2 \alpha_0 \phi'(0)$     sufficient decrease condition

Then stop. $\quad (c \approx 10^{-4})$

② Otherwise, make quadratic approximation to $\phi(\alpha)$ and let $\alpha_1$ be the minimizer of $\phi_q$.     Called $\phi_q$

If $\alpha_1$ satisfies Ⓐ, stop.

③ Otherwise, make cubic approximation to $\phi$, called $\phi_c$, and let $\alpha_2$ minimize $\phi_c$.

If $\alpha_2$ satisfies Ⓐ, stop

else repeat ③.

Cubic functions are good for approximating fens with much curvature.

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

~~Stuff~~ Simulated annealing → more later

# Simplex Method

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

| $f$ only Knowledge | partial Method # | |
|---|---|---|
| $f$ only | CCD (simulated annealing, simplex) | sublinear (?) (linear) Slow |
| $f'$ only | steepest descent | linear |
| partial $f''$ | Quasi-Newton Fisher Scoring | super linear |
| Full $f''$ | Newton | quadratic fast |

~~~~~~~~~~~~~~~~~~~~~~~~~

~~The EM algorithm~~

Fisher Scoring — Poisson regression

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

$\eta_i = \log \mu_i = X_i^T \beta$ ~ $\text{Poisson}(\exp(X^T\beta))$

$\eta = X^T\beta$

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

① Start with $\hat{\mu}_0$

② Set $\hat{z} = \hat{\eta} + (y - \hat{\mu})/\hat{\mu}$

(adjusted response, working response)

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

$\eta = \log \mu$

$\dfrac{d\eta}{d\mu} = \dfrac{1}{\mu}$