hw 2b)   requires   $(A \otimes B)(C \otimes D) = AC \otimes BD$

should prove this identity on the hw

back to optimization...

  A algorithms for finding minima of
  $f(x)$    $f: R^n \to R$

structure of most algorithms goes like this:

  start w/ some $x_0$

  for $k = 0, 1, 2, \ldots$ until convergence

  — Compute $g_k = \nabla f(\vec{x}_k)$

Newton's method
$\left\{ \begin{array}{l} \text{Maybe compute} \\ \text{the hesian (newton's method)} \\ H_k = \nabla^2 f(\vec{x}_k) \text{ & solve} \\ H_k d_k = -g_k \text{ for } d_k \end{array} \right.$

deriv. of f wrt each component of f. ie, x is a vector

$\underline{OR}$

steepest descent
$\left\{ \begin{array}{l} d_k = -g_k \end{array} \right.$

$x_{k+1} = x_k + \alpha_k d_k$ , for some $\alpha_k$ T.B.D.

how to determine "convergence"?

also, note that w/ hesion, (newton's method) need
to check if $g_k^T d_k < 0 \Rightarrow$ we indeed have a
descent direction. (it's guaranteed / by def?
for steepest descent).

If we know we're working w/ a descent direction,
then for a reasonable choice of alpha we're
going to go down.

Convergence test : $\|g_k\|$ is "small enough".
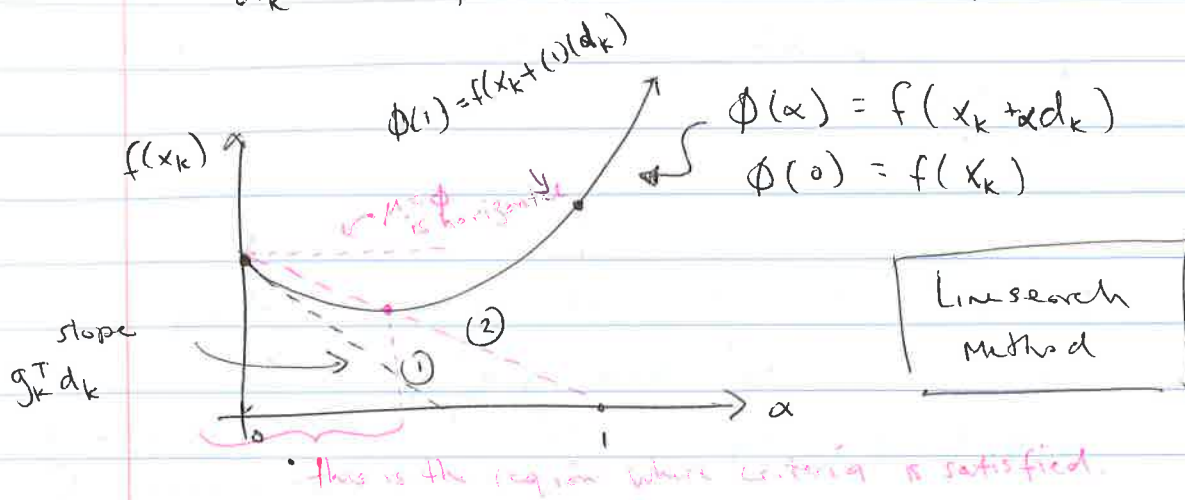small enough might be abs. or relative error.
ie. $\dfrac{\|g_k\|}{\|f(x_k)\|} \leq \hat{\tau}$, tolerance, which depends
on the problem / context.

Cost - assume we have an analytic formula that
includes expressions for $\nabla f$ & $\nabla^2 f$.
there are n iterations for steepest descent,
$n^2$ for newton (? i think this is what he
said...).

So, how to choose $\alpha_k$? (assume we have a $d_k$ already - by whatever method).



$\phi(1) = f(x_k + (1)(d_k))$

$\phi(\alpha) = f(x_k + \alpha d_k)$

$\phi(0) = f(x_k)$

$f(x_k)$

M is horizontal

slope $g_k^T d_k$

① ②

$\alpha$

Line search Method

• This is the region where criteria is satisfied.

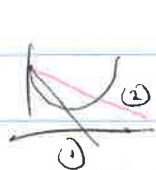In this example $\phi(1) > f(x_k)$, so we want to pick another version of $\alpha$. Strategy?

Taylor Series:

$$\phi(\alpha) = f(x_k + \alpha d_k) = f(x_k) + \alpha g_k^T d_k + \theta(d^2)$$

ignore

Line through $(0, f(x_k))$ w/ slope $g_k^T d_k$ $(< 0)$

Since we know $d_k < 0$, since Line ① is below function $f$ & tangent to $f$ @ $0$, then if we create a new line w/ a larger (but still -'ve) slope, then it is guaranteed

to be above $f^* f$ for at least some period of time. (even if you had something like where $f$ goes below line ① ).



the pink dashed line is a line through $(0, f(x_k))$ w/ slope $\mu g_k^T d_k$ , for $\mu \in (0,1)$. eg. $\mu = \frac{1}{2}$.

| $\mu = 0 \Rightarrow$ horizontal
| $\mu = 1 \Rightarrow$ slope is $g_k^T d_k$, same
|  as original line.

$\mu \approx \frac{1}{2}$ convention.
people research this.
(so, fun!)

Req'nt on $\alpha_k$ :

eq ①

$$f(x_k + \alpha_k d_k) \leq f(x_k) + \mu \alpha_k g_k^T d_k$$

$\underbrace{\qquad\qquad}$ value of $f^*$ line at $\alpha_k$

$\underbrace{\qquad\qquad}$ value of pink line at $\alpha_k$.

"simple-minded approach" $\begin{cases} \text{try } \alpha_k = 1 . & \text{If req't is satisfied, then} \\ \qquad \text{accept.} \\ \text{otherwise, try } \alpha_k \leftarrow \alpha_k / 2 \quad (\text{reduce } \alpha_k) \\ \text{iterate...} \end{cases}$

✗ note: $\exists$ still more book keeping to do, to ensure we have a global min.

For n-dimensions, think of the 2D graph as a cross-section in the direction of $\vec{d_k}$ (vector).

Typical convergence th'm

For a descent method applied to a function f that is bounded below, for which gradient g is Lipschitz continuous, then $\|g_k\| \to 0$

↳ Lipschitz continuous
$$\equiv \|g(x) - g(y)\| \leq L \| x - y \|$$
for some $L > 0$
$\forall$ x,y in the region of interest.

∃ better ways to choose $\alpha$? viz. simple minded approach discussed earlier — want a diff way that could be smarter ... ie, if we can use our knowledge of f to reduce stupid guesses, we are happier.

In particular, we'll have to compute eq² ① less often.

Cost of computing ① ?   Well,

$f(x_k)$ - free

$\max_{\alpha_k} g_k^T d_k$ - cheap

but   $f(x_k + \alpha_k d_k)$   (LHS) might be
arbitrarily complex $f^{\circ} f$.

So , an alternative approach to meet the req't
of eq? ① :

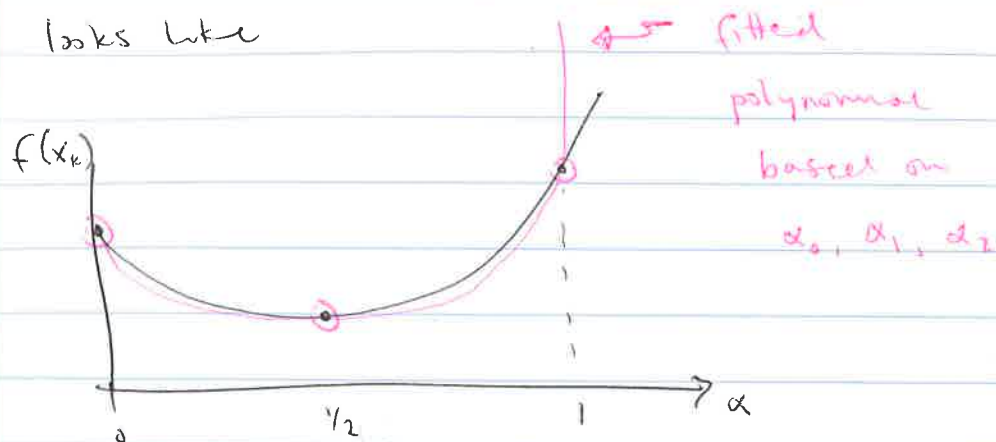start w/     $\alpha_{k_1} = 1$

$\alpha_{k_2} = $ something like
$\frac{1}{2} \alpha_{k_1}$

Now, we have candidate values:

$\alpha_k = 0$   $f(x_k)$         $\alpha_k = 1$   $f(x_k + d_k)$
$\alpha = \frac{1}{2}$   $f(x_k + \frac{1}{2} d_k)$

Let $p(\alpha) = $ quadratic polynomial that
interpolates 3 points. Then, find $\alpha$
that minimizes $p(\alpha)$.   call this
$\alpha_{k_3}$.

looks like



the fitted polynomial makes use of values
we've already computed.