

## CS205 Review Session #4 Notes

### Secant Formula

Consider the following alternate form of the secant method:

$$x_{k+1} = \frac{x_{k-1}f(x_k) - x_k f(x_{k-1})}{f(x_k) - f(x_{k-1})}$$

This formula leads to an indefinite form ( $\frac{0}{0}$ ) as the current iterate asymptotically approaches the true solution. If  $x_*$  is an analytic root of  $f(x)$  s.t.  $f(x_*) = 0$ , we can see that  $x_k \approx x_{k-1} \approx x_*$  and  $f(x_k) \approx f(x_{k-1}) \approx 0$  as  $x_{k-1}, x_k \rightarrow x_*$ .

This observation is significant, since the aforementioned update formula relies on the division of two quantities that are very close to zero, which may lead to numerical instability. Similar behavior may be seen in the more usual form of the secant update formula:

$$x_{k+1} = x_k - f(x_k) \frac{x_k - x_{k-1}}{f(x_k) - f(x_{k-1})}$$

### Mean Value Theorem

Consider a real-valued function  $f(x)$  that is differentiable and continuous on the interval  $[x_k, x_{k-1}]$ . Then, the **Mean Value Theorem** states that:

$$\exists \hat{x} \in [x_k, x_{k-1}] \text{ s.t. } f(x_k) - f(x_{k-1}) = f'(\hat{x})(x_k - x_{k-1})$$

Furthermore, it is clearly the case that  $f'(\hat{x}) \approx f'(x_*)$  as  $x_k, x_{k-1} \rightarrow x_*$ .

### Fixed-Point Iteration

One nice property of root finding problems is that they can always be reformulated in terms of an equivalent **fixed-point** problem. In particular, if we wish to find  $x$  such that  $f(x) = 0$ , we can pick an appropriate function  $g$  such that  $g(x) = x \Leftrightarrow f(x) = 0$ . In fact, there are infinitely many  $g$  functions we can pick:

$$\begin{aligned} g(x) &= x + f(x) \\ g(x) &= x(1 + f(x)) \\ g(x) &= xe^{f(x)} \\ g(x) &= x + f(x)f'(x) \end{aligned}$$

The fixed-point formulation may, under certain conditions, be quite convenient. In particular, we may be able to use **fixed-point iteration** to find the solution to both problems. Given some initial guess  $x_0$  that's sufficiently close to the exact solution  $x_*$ , we iterate as follows:

$$x_k = g(x_{k-1})$$

To see that this process will converge to the correct solution, consider the error at the  $k^{\text{th}}$  iteration:

$$\begin{aligned} e_k &= x_k - x_* \\ &= g(x_{k-1}) - g(x_*) \end{aligned}$$

Now, by the Mean Value Theorem, we have (for some  $\hat{x}_k \in [x_{k-1}, x_*]$ ):

$$\begin{aligned} g(x_{k-1}) - g(x_*) &= g'(\hat{x}_k)(x_{k-1} - x_*) \\ &= g'(\hat{x}_k)e_{k-1} \end{aligned}$$

We may now see that it suffices to have  $|g'(x_*)| < 1$ . If this is the case, and we pick  $x_0$  sufficiently close to  $x_*$ , then there exists some constant  $c$  such that  $g'(\hat{x}_k) \leq c < 1$  for  $k \in \{0, 1, \dots\}$ . From the above, this gives us:

$$|e_k| = g'(\hat{x}_k) |e_{k-1}| \leq c |e_{k-1}| \leq \dots \leq c^k |e_0|$$

This will clearly converge if  $c < 1$ .

## Optimization Criteria

Some possible criteria used to characterize the effect of an optimization step are given below:

1. Does the current guess for  $x$  approximate the optimal  $x_{\min}$  better than the previous one? That is  $|x_{k+1} - x_{\min}| \stackrel{?}{\gtrless} |x_k - x_{\min}|$
2. Does the current value  $f(x)$  approximate the optimal  $f(x_{\min})$  better than the previous one? That is  $|f(x_{k+1}) - f(x_{\min})| \stackrel{?}{\gtrless} |f(x_k) - f(x_{\min})|$
3. Is the size of the update in  $x$  decreasing? That is  $|x_{k+1} - x_k| \stackrel{?}{\gtrless} |x_k - x_{k-1}|$
4. Is the relative improvement in  $f(x)$  decreasing? That is  $|f(x_{k+1}) - f(x_k)| \stackrel{?}{\gtrless} |f(x_k) - f(x_{k-1})|$

## The Metric Tensor

In the derivation of the Conjugate Gradient method, we frequently encounter expressions of the form  $\mathbf{x}^T \mathbf{A} \mathbf{y} = \mathbf{x} \cdot \mathbf{A} \mathbf{y}$  where  $\mathbf{x}$  and  $\mathbf{y}$  are vectors and  $\mathbf{A}$  is the symmetric positive definite matrix representing the coefficients of the linear equations we wish to solve. In particular, we make extensive use of the concept of **A-orthogonality**: two vectors  $\mathbf{x}$  and  $\mathbf{y}$  are **A-orthogonal** (or **conjugate** with respect to  $\mathbf{A}$ ) if  $\mathbf{x} \cdot \mathbf{A} \mathbf{y} = \mathbf{y} \cdot \mathbf{A} \mathbf{x} = 0$ .

Recall that the fundamental notion of **distance** in a vector space can be built from the definition of an **inner product**. Given some inner product  $\langle \cdot, \cdot \rangle$ , we can define (in the usual fashion):

$$\langle \mathbf{x}, \mathbf{y} \rangle = |\mathbf{x}| |\mathbf{y}| \cos \theta_{xy}$$

This gives us a straightforward extension to the notion of **length**, since we then have:

$$\langle \mathbf{x}, \mathbf{x} \rangle = |\mathbf{x}| |\mathbf{x}| \cos \theta_{xx} = |\mathbf{x}|^2$$

This, in turn, allows us to define the **distance** along a parameterized curve  $\mathbf{p}(t)$ :

$$L = \int_a^b |\mathbf{p}'(t)| dt$$

Any inner product must satisfy, in general, only the following set of properties:

1.  $\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{y}, \mathbf{x} \rangle$
2.  $\langle \mathbf{x}, \mathbf{x} \rangle \geq 0$
3.  $\langle \mathbf{x}, \mathbf{x} \rangle = 0 \iff \mathbf{x} = \mathbf{0}$
4.  $\langle c\mathbf{x}, \mathbf{y} \rangle = c\langle \mathbf{x}, \mathbf{y} \rangle$
5.  $\langle \mathbf{x} + \mathbf{z}, \mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{y} \rangle + \langle \mathbf{z}, \mathbf{y} \rangle$

From this rather loose set of constraints, it should be obvious that infinitely many such inner products exist for any given space. In fact, an inner product is uniquely defined by a **metric tensor**, which (for our purposes) can be thought of simply as an  $n \times n$  matrix  $\mathbf{g}$  that induces an associated inner product by:

$$\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbf{g}} = \mathbf{x}^T \mathbf{g} \mathbf{y}$$

When  $\mathbf{g} = \mathbf{I}$  we have the usual Euclidean dot product, since:

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^T \mathbf{I} \mathbf{y} = \mathbf{x}^T \mathbf{y} = \mathbf{x} \cdot \mathbf{y}$$

From the list of constraints above, we can see that  $\mathbf{g}$  cannot be a completely arbitrary matrix. Constraint (1) implies that  $\mathbf{g}$  must be symmetric, and constraint (2) is the very definition of positive definiteness. With these observations, the concept of conjugacy becomes more intuitively understandable. Two vectors are conjugate with respect to a symmetric positive definite matrix  $\mathbf{A}$  precisely if they are orthogonal under the inner product induced by  $\mathbf{A}$ .