

## Chapter 2

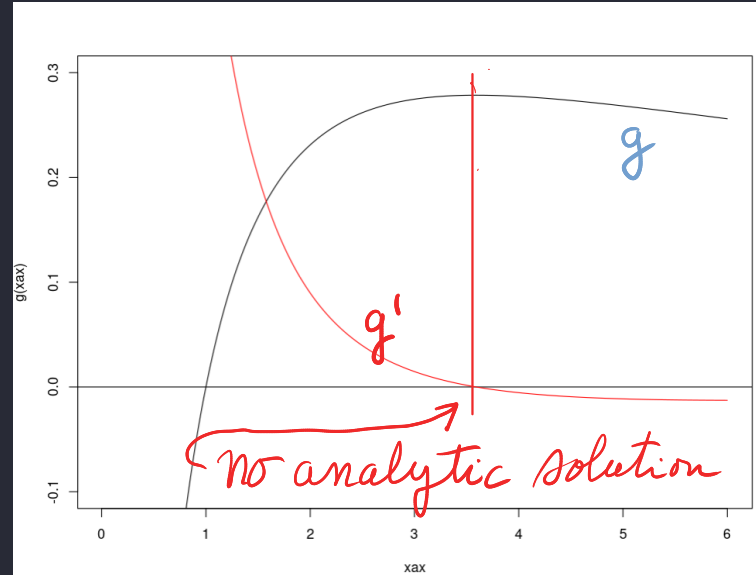
- Optimizing  $f: \mathbb{R}^P \rightarrow \mathbb{R}$   
where  $f$  smooth - differentiable
- Frequent application:
  - Optimize log-likelihood to get MLE.
  - e.g. Find solution,  $\hat{\theta}$ , to  
the score equation  $l'(\hat{\theta}) = \underline{0}$  } root problem

# Univariate $\theta$

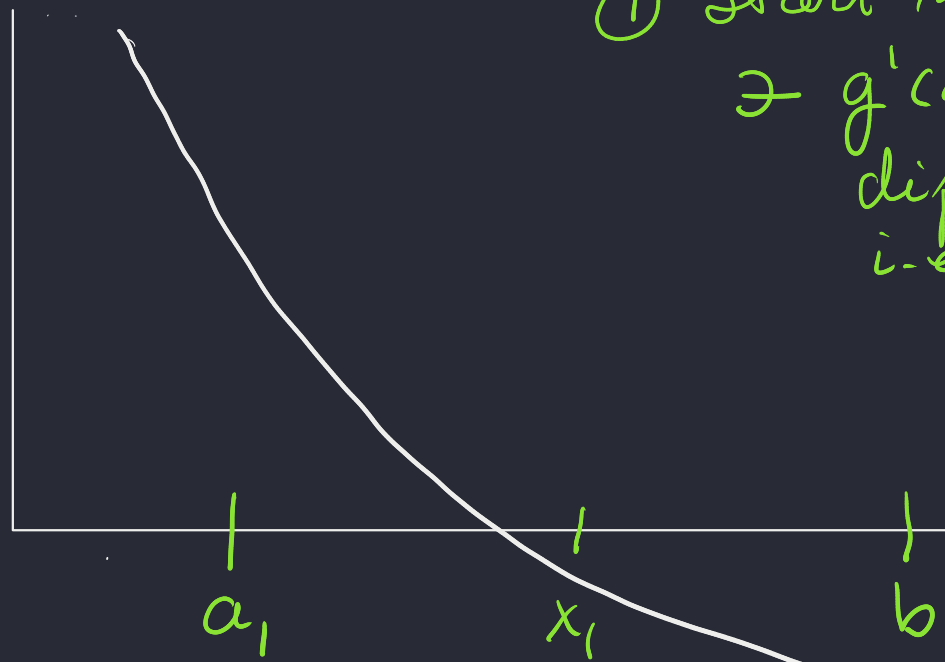
Example:

$$g(x) = \frac{\log(x)}{1+x}$$

$$g'(x) = \frac{\frac{1}{x}(1+x) - \log(x)}{(1+x)^2}$$



# Bisection method : Use $g'$



① Start with  $a_1, b_1$   
 $\exists g'(a_1) \neq g'(b_1)$  have  
different signs  
i.e.  $g'(a_1) \cdot g'(b_1) < 0$

② Bisect  
 $x_1 = (a_1 + b_1) / 2$

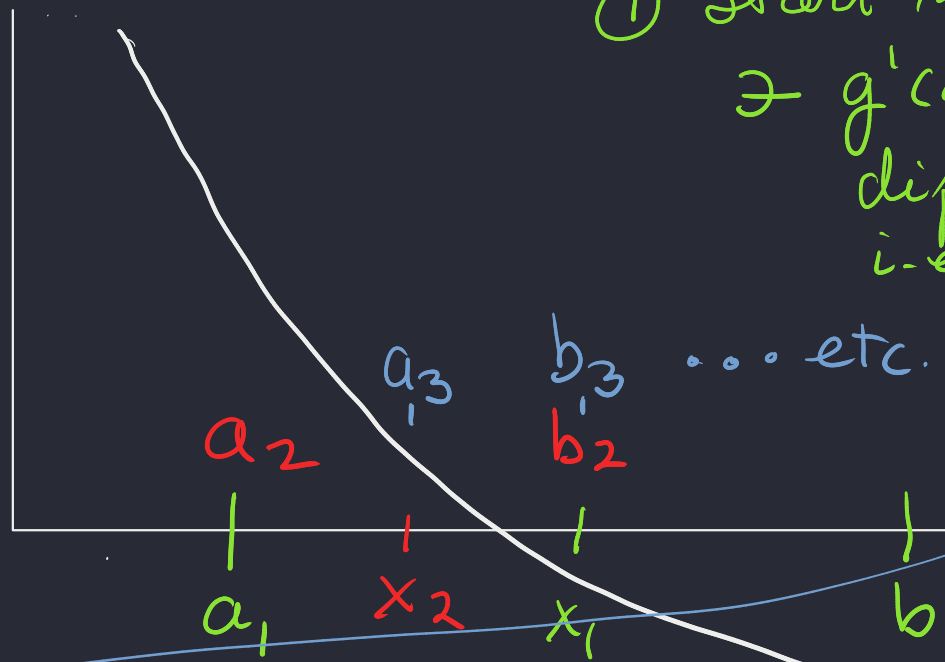
③ Check  
 $g(x_1) \cdot g(b_1)$

If  $< 0$  let  $a_2 = x_1, b_2 = b_1$

if  $> 0$  let  $a_2 = a_1, b_2 = x_1$

Repeat.

# Bisection method : Use $g'$



① Start with  $a_1, b_1$   
 $\exists g'(a_1) \neq g'(b_1)$  have  
different signs  
i.e.  $g'(a_1) \cdot g'(b_1) < 0$

② Bisect  
 $x_1 = (a_1 + b_1) / 2$

③ Check  
 $g(x_1) \cdot g(b_1)$

If  $< 0$  let  $a_2 = x_1, b_2 = b_1$   
if  $> 0$  let  $a_2 = a_1, b_2 = x_1$

Repeat.

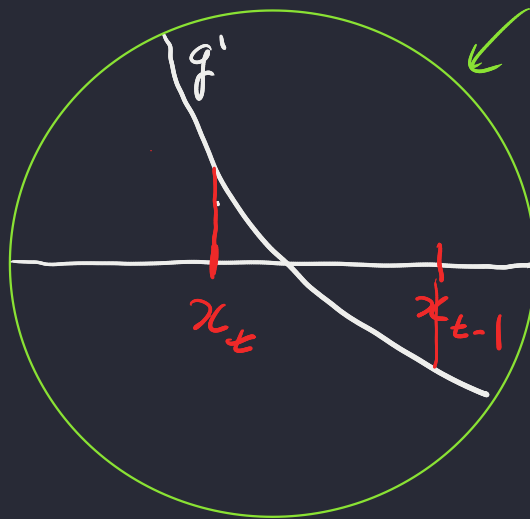
But code :

$$x_1 = a_1 + \frac{b_1 - a_1}{2}$$

# Convergence Criterion:

$$|x_t - x_{t-1}| \text{ small}$$

$$|g'(x_t)| \text{ small}$$



microscope

"Absolute" convergence criterion  $|x_t - x_{t-1}| < \epsilon$

- Choose small  $\epsilon$

- Converge if  $|x_t - x_{t-1}| < \epsilon$

- Bisection must converge since

$$|x_t - x_{t-1}| < 2^{-t} |b_1 - a_1|$$

# "Relative Convergence"

$$a) \frac{|x_t - x_{t-1}|}{|x_t|} < \epsilon$$

— Independent of units of units

— Not good is solution at 0!

$$b) \frac{|x_t - x_{t-1}|}{|x_t| + \epsilon} < \epsilon$$

Caveats: - Of many roots, will find one  
- not necessarily max  
-  $g'$  must be continuous

# Newton's method

Find  $x_*$  such that  
 $g(x_*)$  is max of  $g(x)$

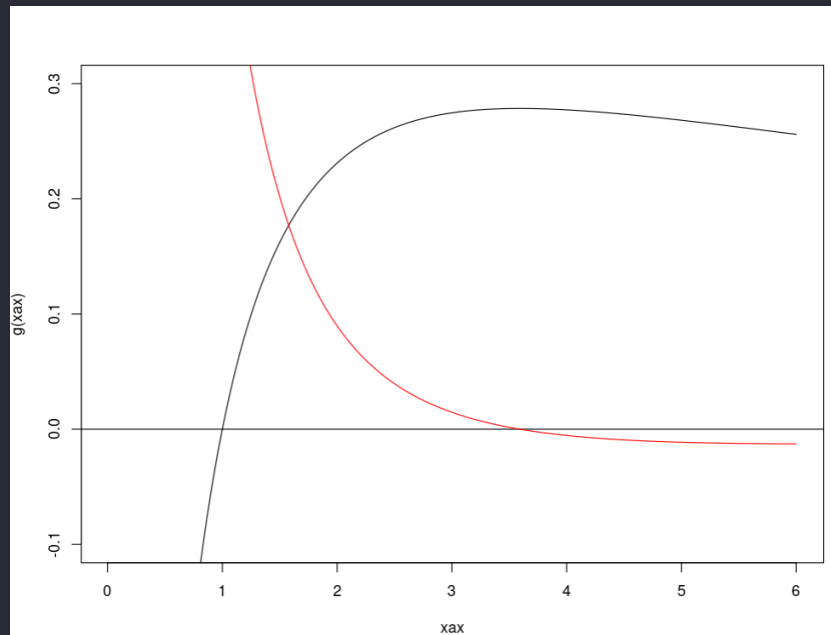
i.e.  $x_* = \operatorname{argmax}(g)$

OR Find  $x_*$  such that

$$g'(x_*) = 0$$

i.e. 1) max of  $g$  by using quad. approx of  $g$

or 2) Root of  $g'$  by using a linear approx of  $g'$





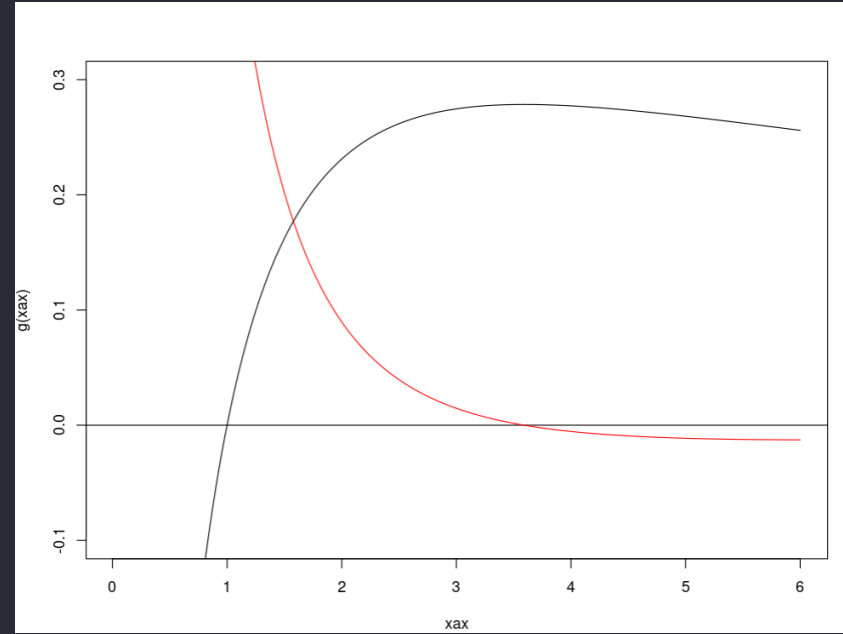
Start somewhere

e.g.  $x_1 = 2$

1) Do a quadratic approximation of  $g$  at  $x_1$

2) Solve for max using quad. approx.

OR: 1) Do a linear approximation of  $g'$  at  $x_1$   
2) Solve for  $\hat{g}'(x) = 0$



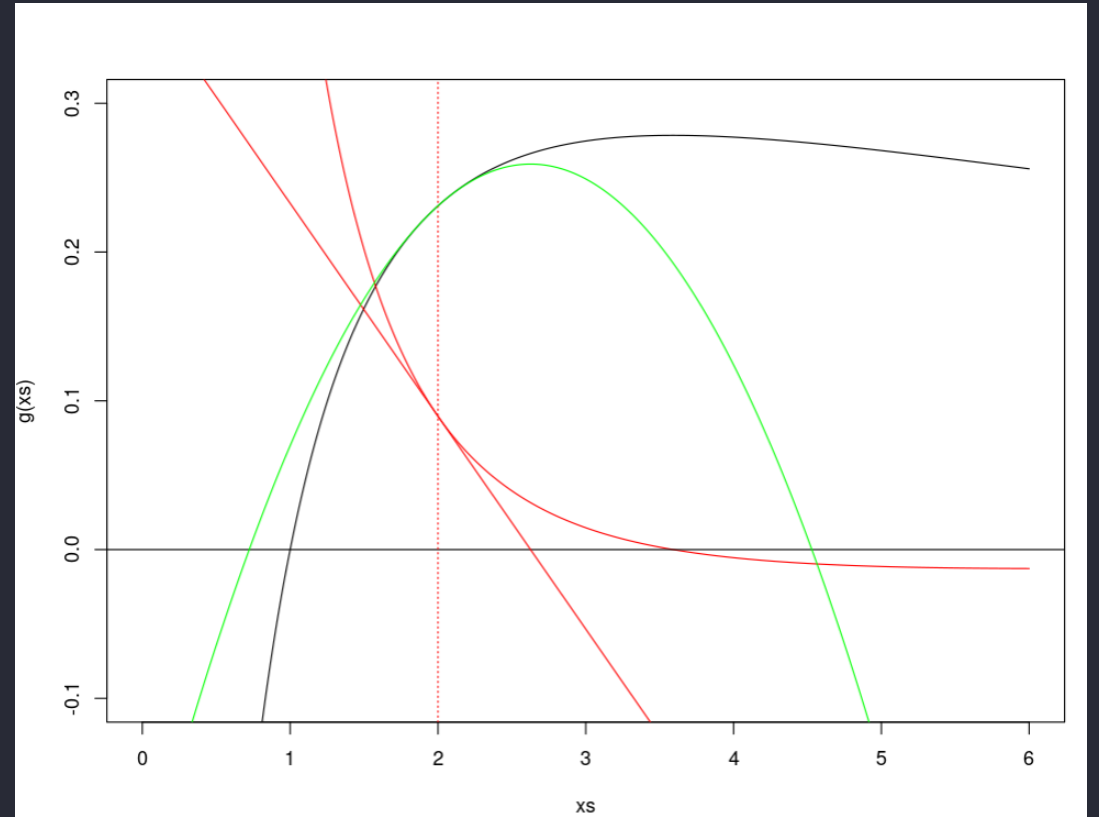
# Quadratic approx

$$\begin{aligned}\hat{g}(x) &= g(x_1) \\ &+ g'(x_1)x(x-x_1) \\ &+ g''(x_1) \cdot (x-x_1)^2 / 2\end{aligned}$$

Quad:  $a + bx + cx^2$   
if  $c < 0$  then max at  
 $x = -b/2c$

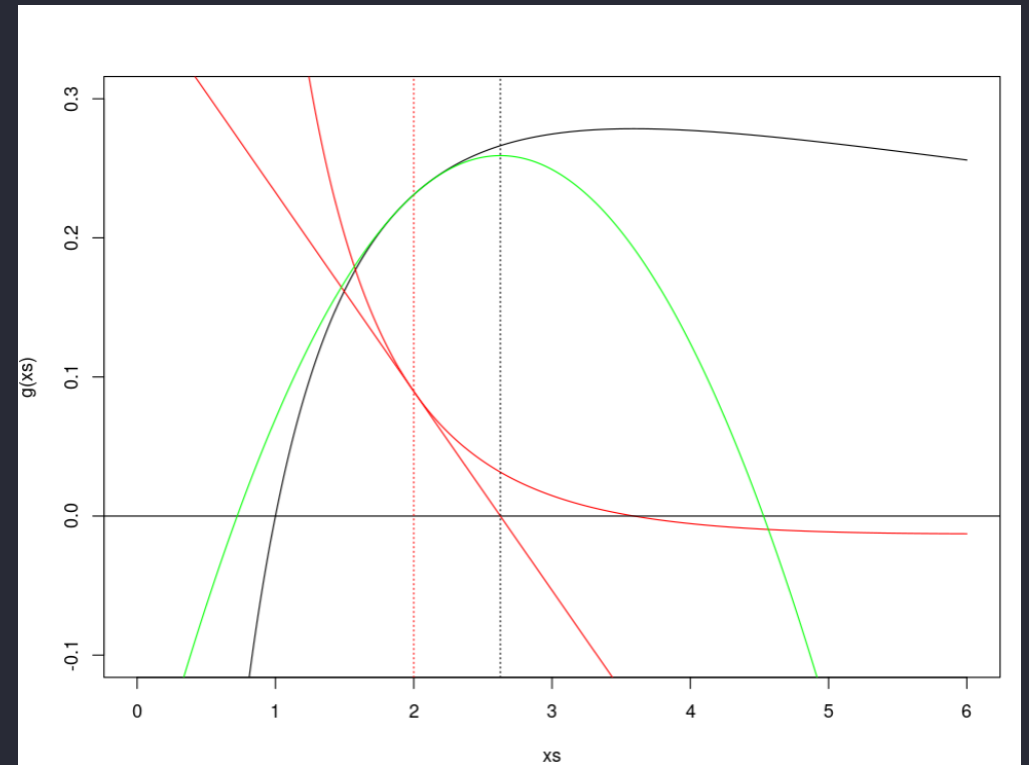
So max of  $\hat{g}$  is at  $(x-x_1) = -g'(x_1)/g''(x_1)$

i.e.  $x_2 = x_1 - g'(x_1)/g''(x_1)$

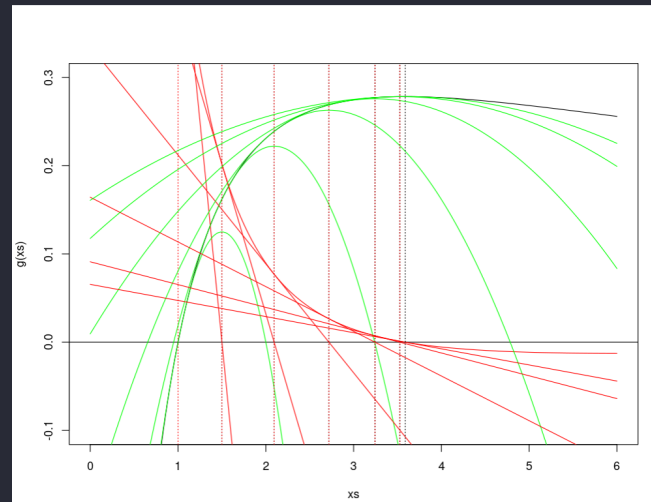
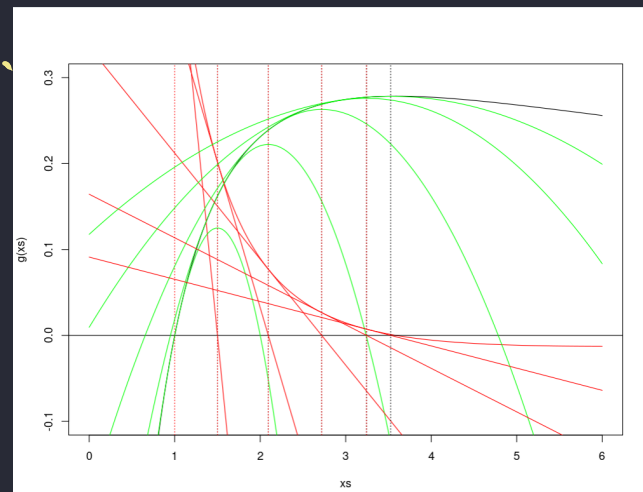
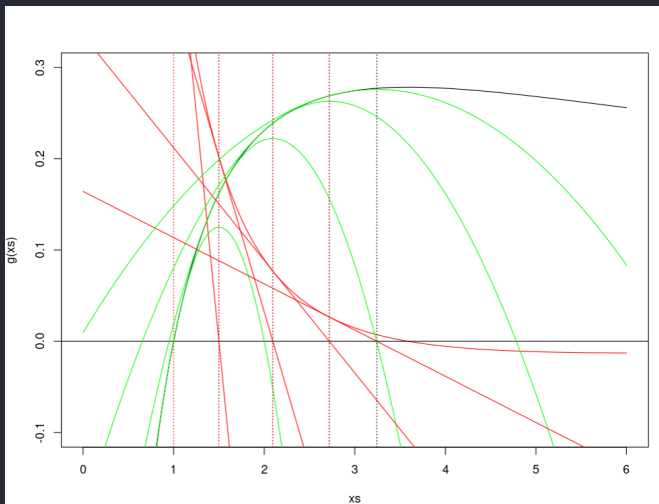
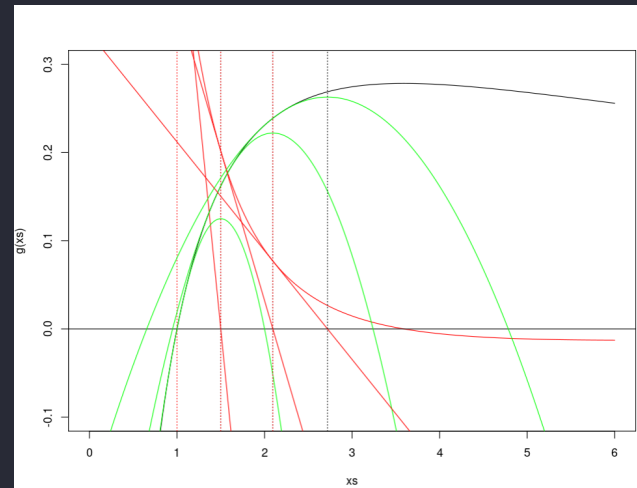
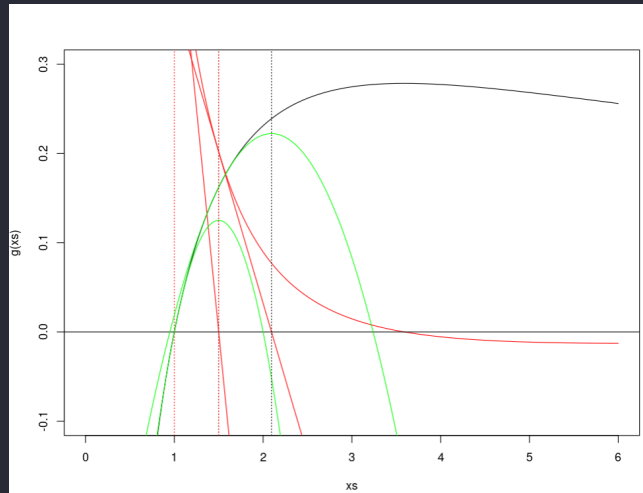
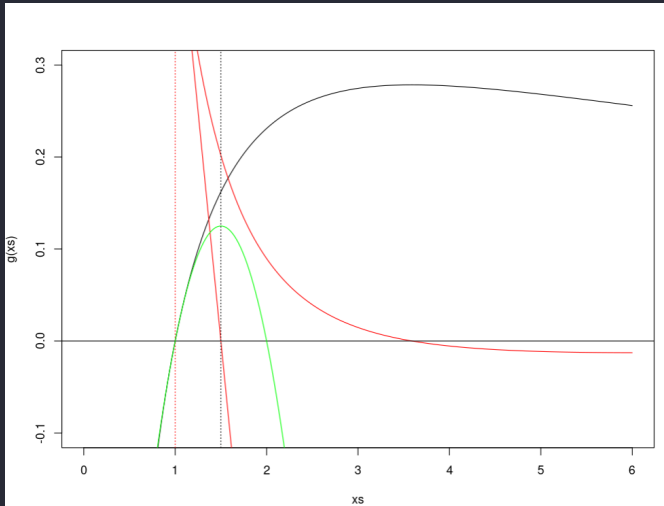


EXER: Show solving for  $g'(x) = 0$   
yields the same solution.

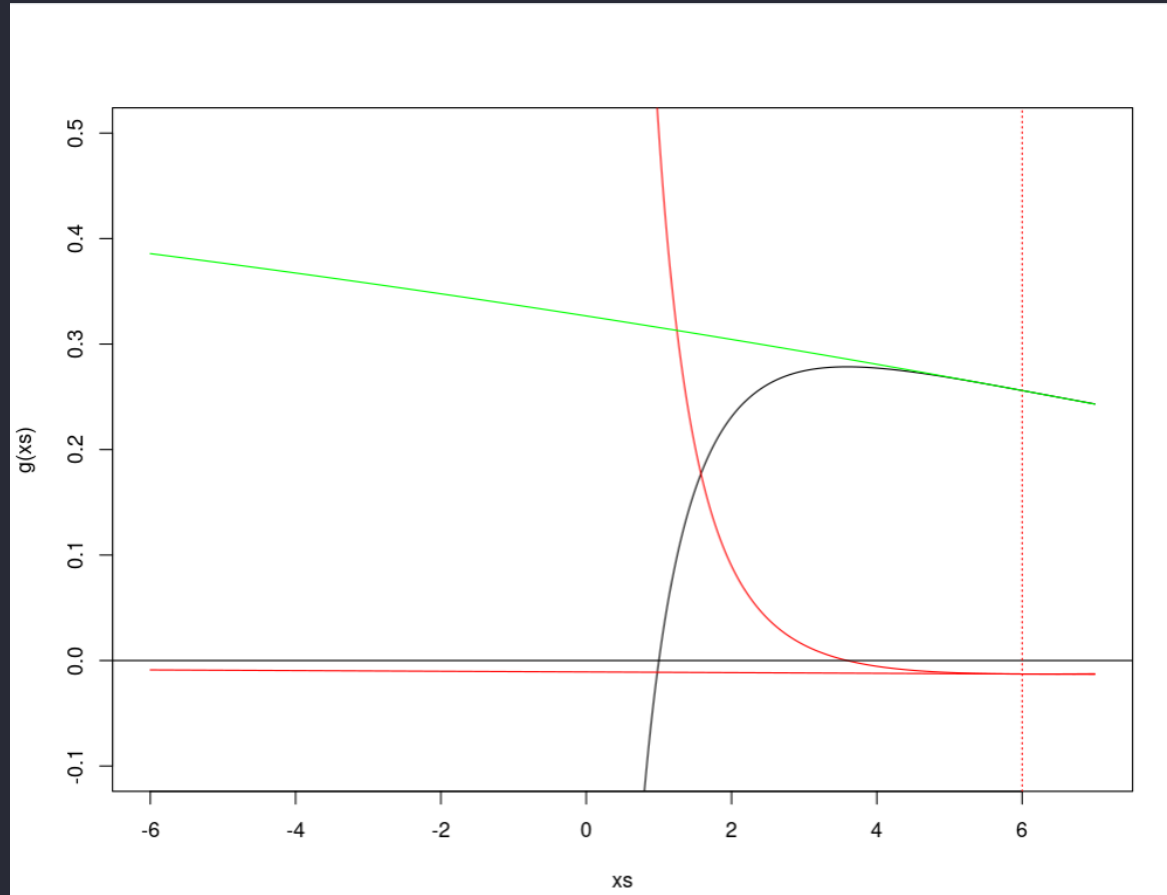
Repeat



Starting at  $x = 1$



What happens if we start at  $x = 6$



6 looked like a perfectly reasonable starting value.

Generalize to multivariate  $x$

$$g(x_1, x_2)$$

$$g'(x_1, x_2) = \begin{pmatrix} \frac{\partial g}{\partial x_1} & \frac{\partial g}{\partial x_2} \end{pmatrix}$$

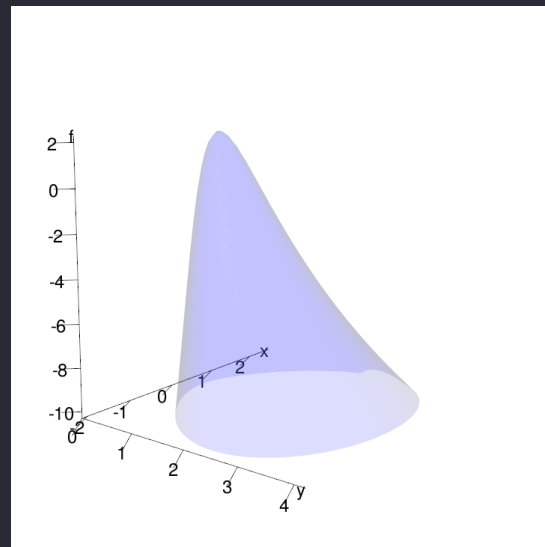
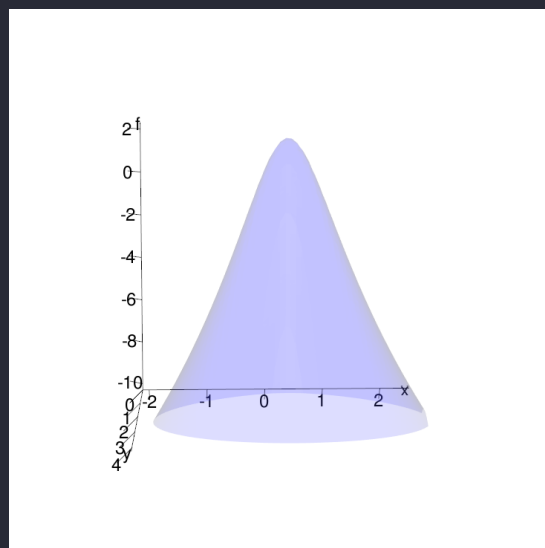
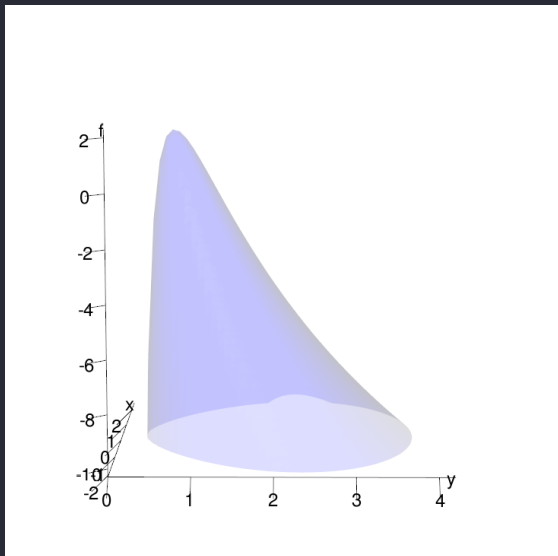
$$g''(x_1, x_2) = \begin{pmatrix} \frac{\partial^2 g}{\partial x_1^2} & \frac{\partial^2 g}{\partial x_1 \partial x_2} \\ \frac{\partial^2 g}{\partial x_1 \partial x_2} & \frac{\partial^2 g}{\partial x_2^2} \end{pmatrix}$$

## Quad approximation of $\tilde{x}_0$

$$\begin{aligned}\hat{g}_{x_0}(\tilde{x}) &= g(\tilde{x}_0) + \underbrace{g'(\tilde{x}_0)}_{1 \times 2 \text{ vec}} \underbrace{(\tilde{x} - \tilde{x}_0)}_{2 \times 1 \text{ vector}} \\ &\quad + \underbrace{(\tilde{x} - \tilde{x}_0)}_{1 \times 2} \underbrace{g''(\tilde{x}_0)}_{2 \times 2} \underbrace{(\tilde{x} - \tilde{x}_0)}_{2 \times 1} / 2\end{aligned}$$

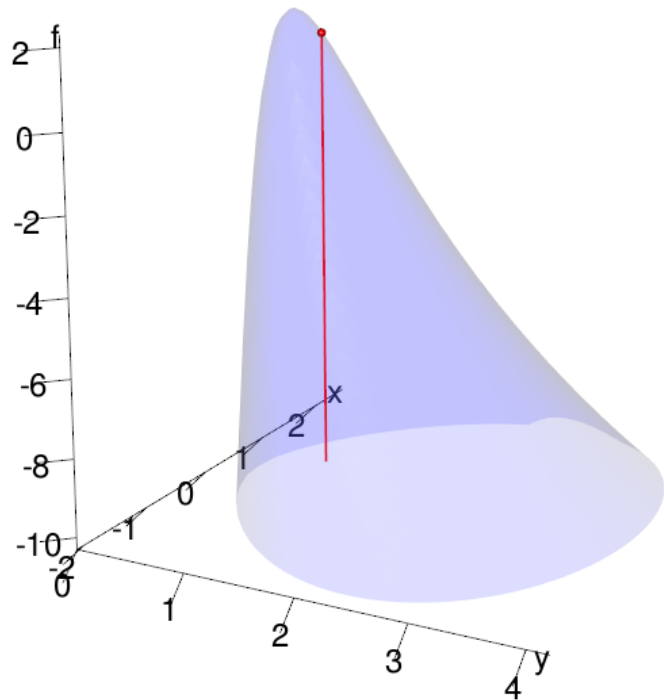
Has a single max at  $\tilde{x}$  if  $-g''(\tilde{x}_0)$  is positive definite

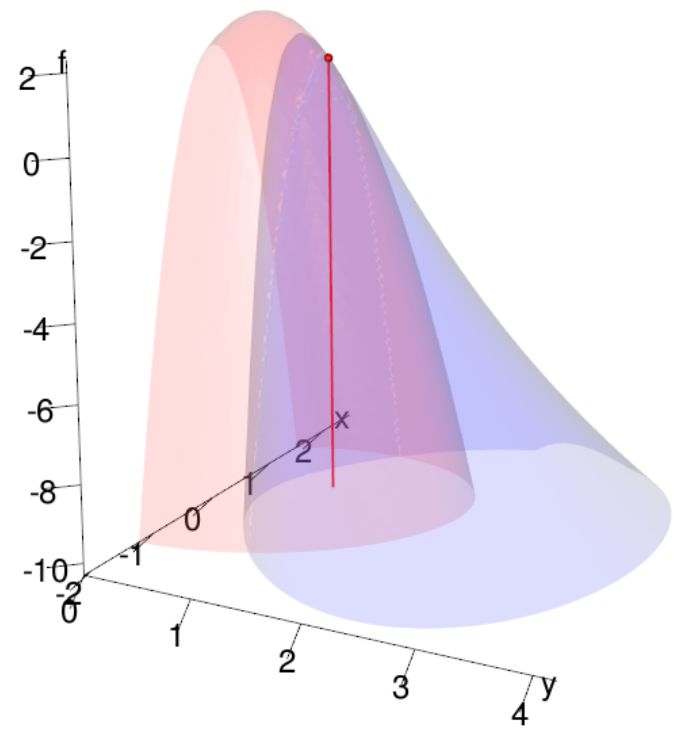
$$\hat{\tilde{x}} = \tilde{x}_0 - [g''(\tilde{x}_0)]^{-1} [g'(\tilde{x}_0)]^T$$



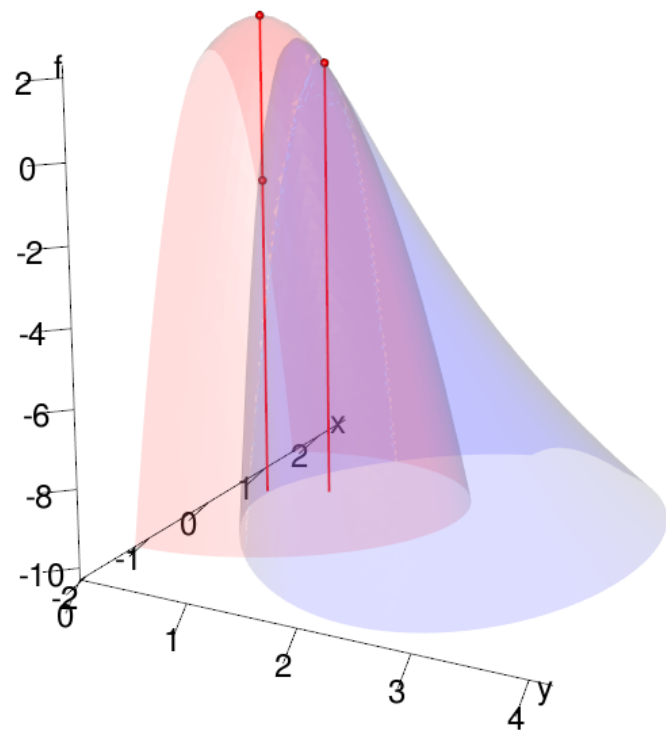


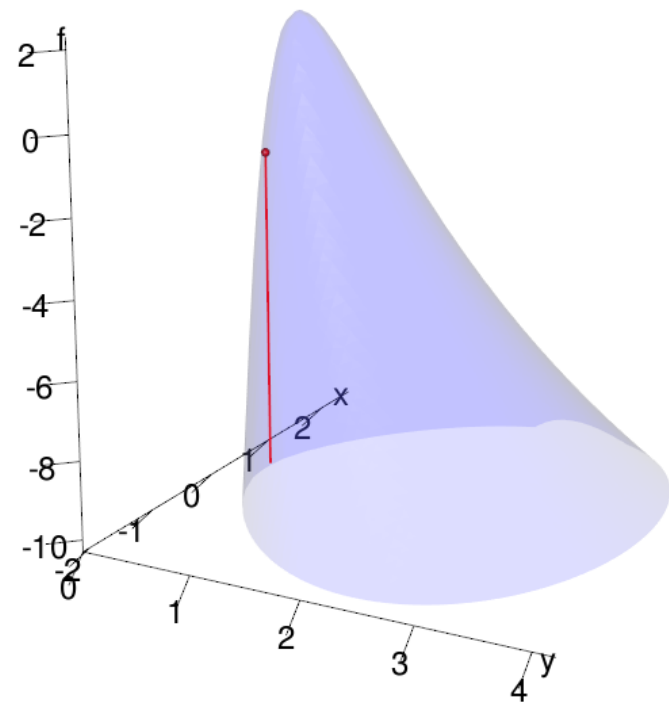
Start at  $\mu_1 = 0.5$ ,  $\sigma_1 = 0.9$

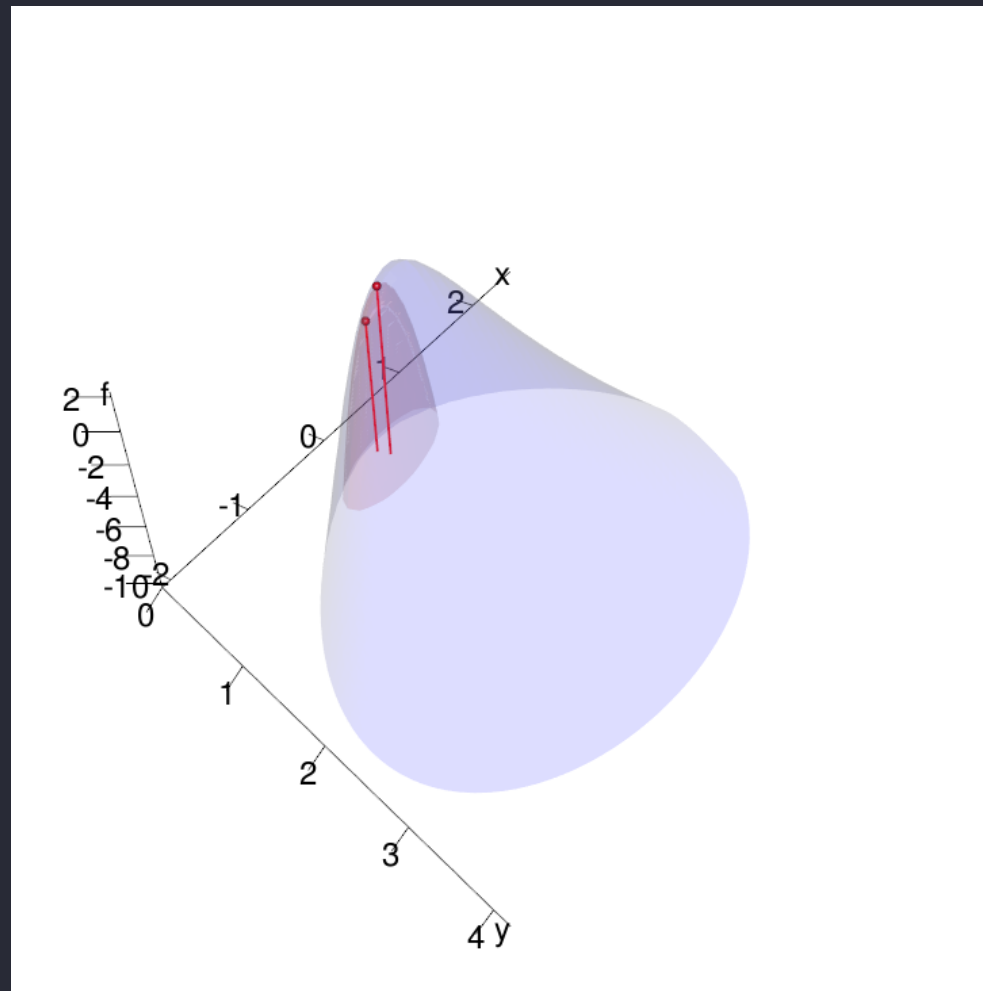
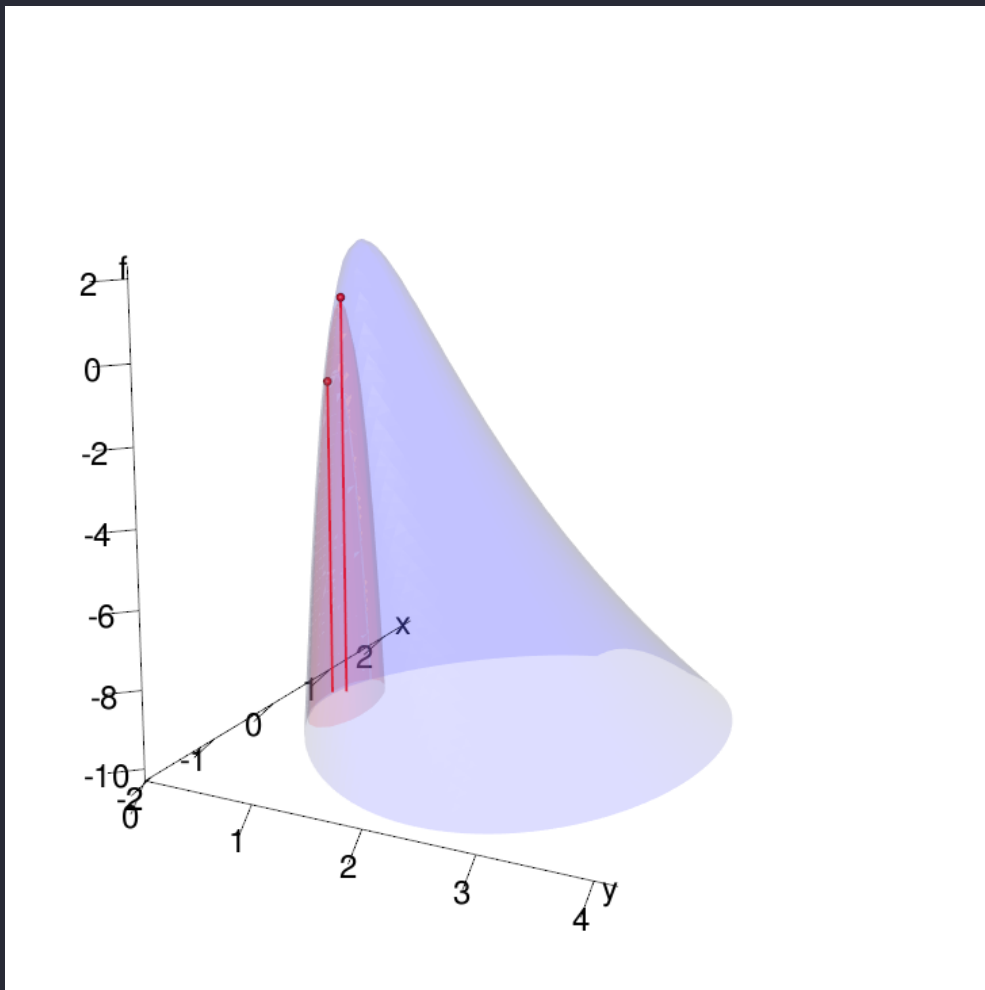




```
> Deriv(ll, x)
      [,1] [,2]
[1,] -1.520725 -4.167008
> drop(Deriv(function(x) Deriv(ll,x), x)) # drops dimensions of size 1
      [,1] [,2]
[1,] -12.345679 3.379392
[2,] 3.379392 -10.801332
> xhat <- x - c( solve(drop(Deriv(function(x) Deriv(ll,x),x))) %*% t(Deriv(ll,x)) )
> xhat
[1] 0.2497914 0.4359312
```

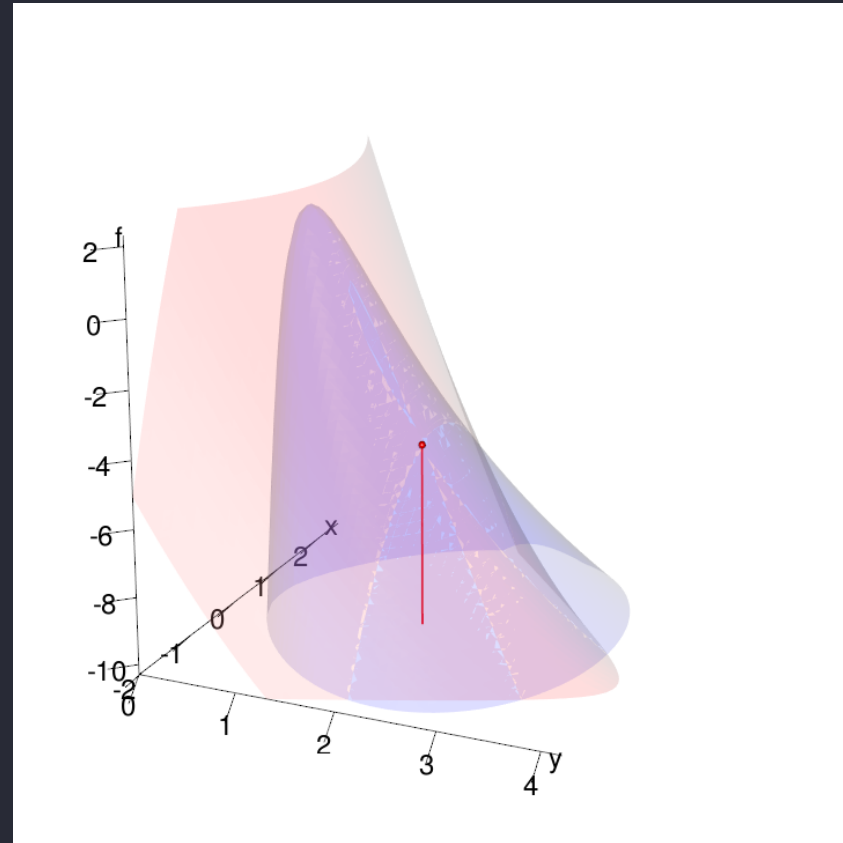
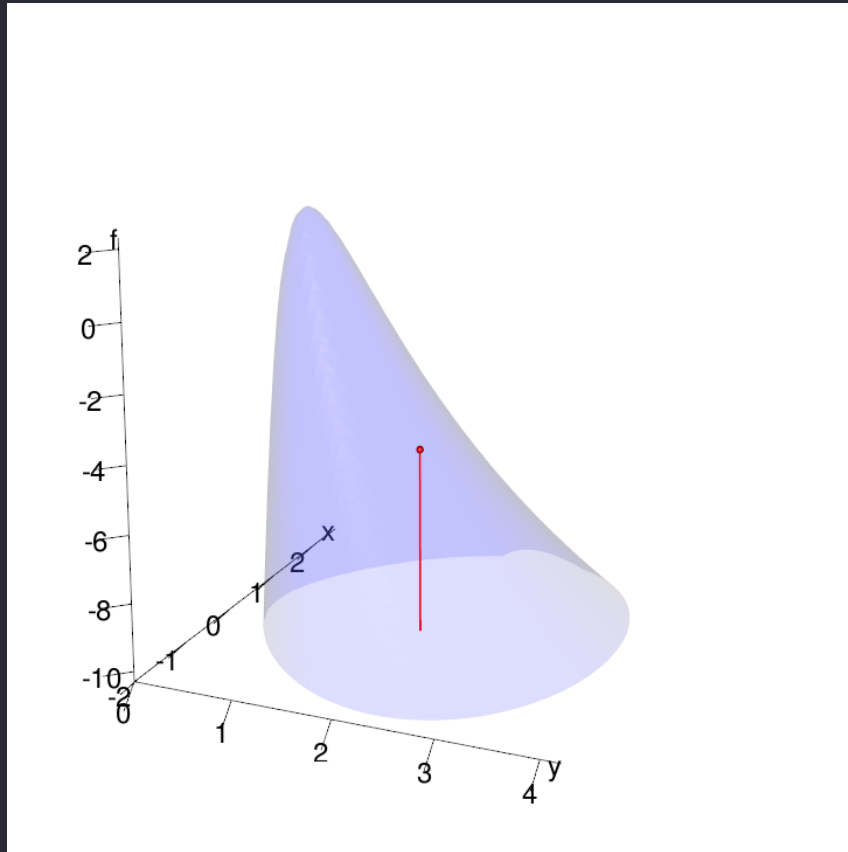






... gets closer to  $\text{argmax}(l)$  ...

What if we start at  $\mu_1 = 0, \sigma_1 = 1$



```
> xhat <- x - c( solve(drop(Deriv(function(x) Deriv(ll,x),x))) %% t(Deriv(ll,x)) )  
> xhat  
[1] -0.6541702  4.7360095
```

o o o ? . . . ?

Newton's method needs a good starting point - close enough to the maximum for the quadratic approximation to

1) have a negative definite Hessian  
positive - for a minimum

2) have a maximum that doesn't take you outside the domain in which further quadratic approximation eventually take you to the maximum.

When will N/A work?



Recall Taylor Series expansion not a approximation

If  $f, f', f''$  are continuous in  $(x_0 \pm \epsilon)$

$\forall x \in (x_0 \pm \epsilon), \exists x_1 \in (x_0 \pm \epsilon)$

$$\exists f(x) = f(x_0) + f'(x_0)(x - x_0) + f''(x_1)(x - x_0)^2 / 2$$

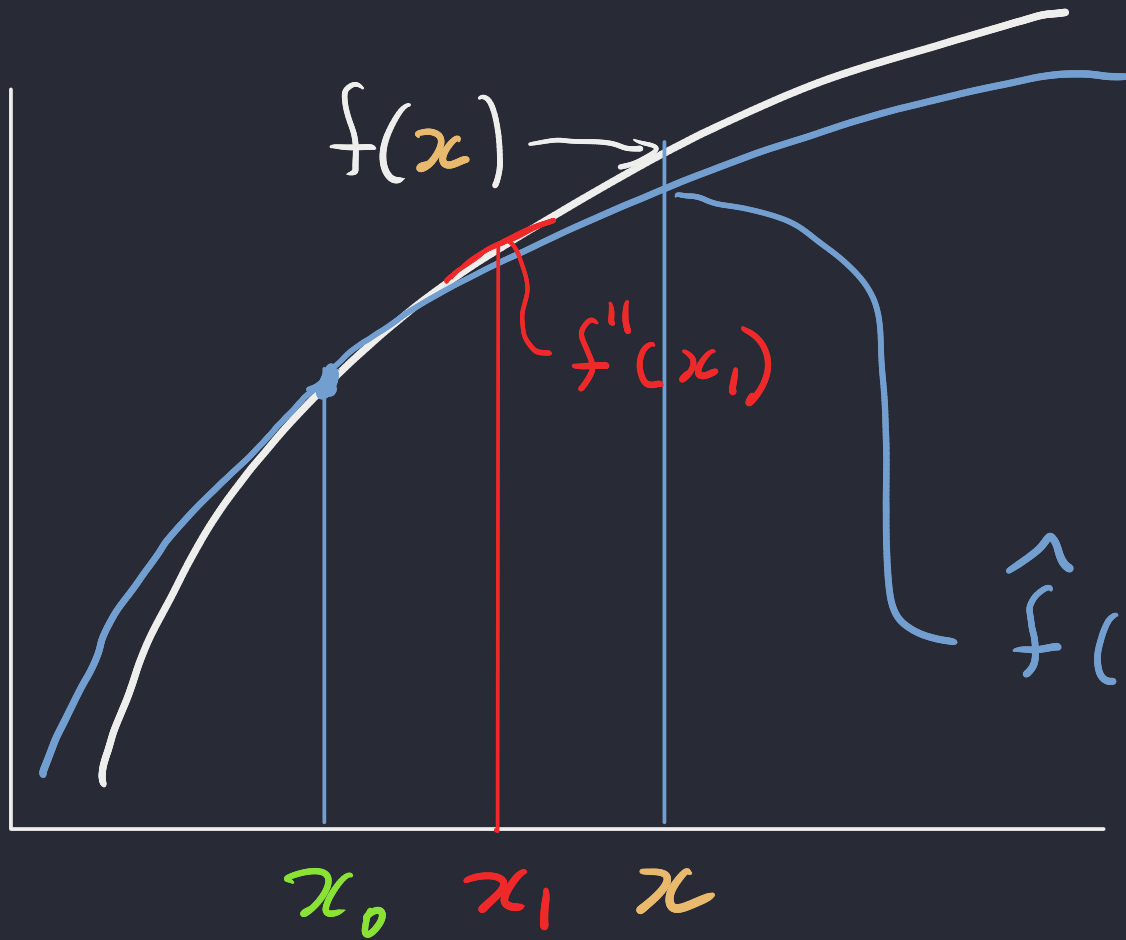
Taylor series expansion

Note:

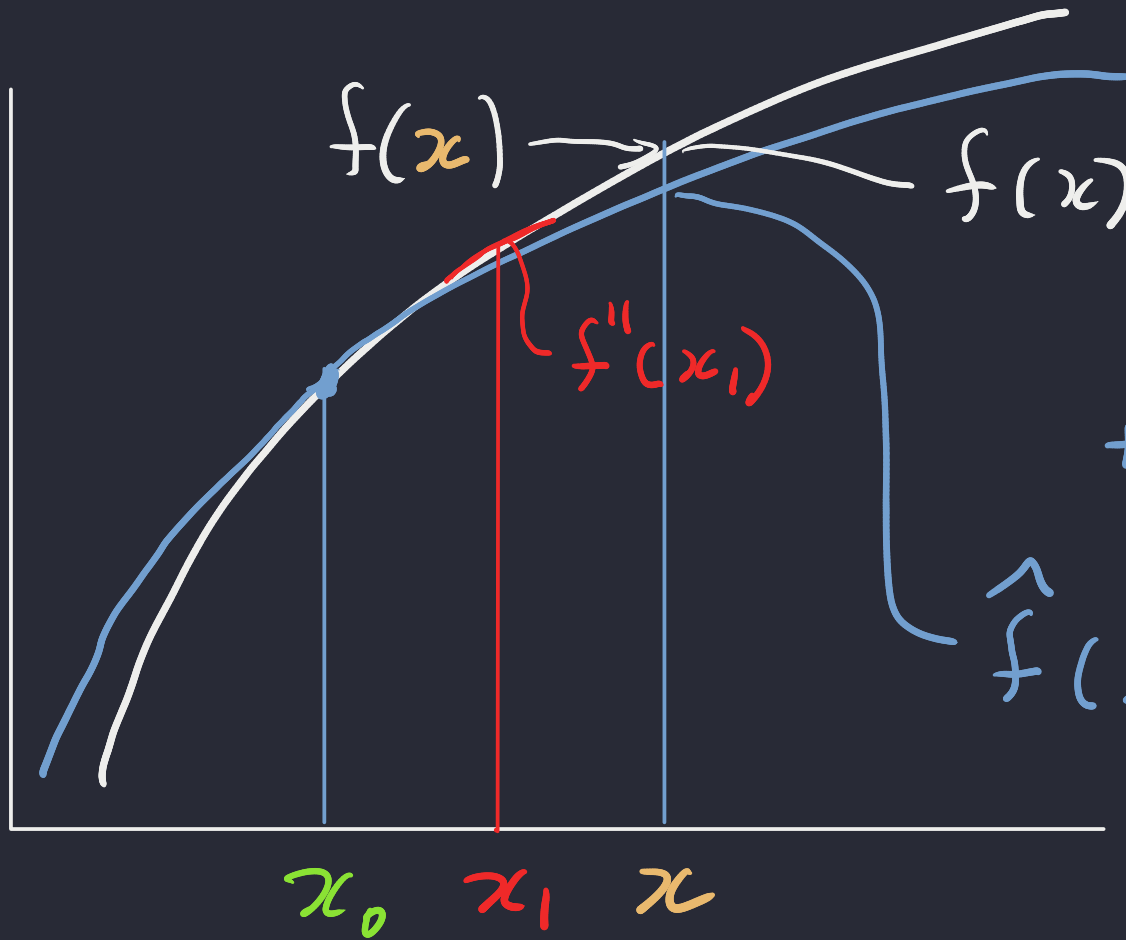
$$\hat{f}(x) = f(x_0) + f'(x_0)(x - x_0) + f''(x_0)(x - x_0)^2 / 2$$

*← compare*

Taylor series approximation



$$\hat{f}(x) = f(x_0) + f'(x_0)(x-x_0) + f''(x_0)(x-x_0)^2/2$$



$$f(x) = f(x_0) + f'(x_0)(x - x_0) + f''(x_1)(x - x_0)^2 / 2$$

$$\hat{f}(x) = f(x_0) + f'(x_0)(x - x_0) + f''(x_0)(x - x_0)^2 / 2$$

## BACK TO:

When can we be sure that  
Newton's method will work?

Let  $x^*$  be  $\operatorname{argmax}(g)$ ,  $g \in C^3$ .

So  $g'(x^*) = 0$ . Suppose  $g''(x^*) < 0$ .

By continuity,  $\exists \varepsilon > 0$

$$\exists \delta > 0 \quad \forall x \in (x^* - \varepsilon, x^* + \varepsilon) \quad g''(x) < -\delta$$

Let  $\varepsilon_t = x_t - x^*$  for  $|\varepsilon_t| < \varepsilon$

Apply Taylor expansion to  $g'$

$$0 = g'(x^*) = g'(x_t) + g''(x_t)(x^* - x_t) + \frac{1}{2}(x^* - x_t)^2 g'''(x^b)$$

for  $x^b$  between  $x^*$  &  $x_t$

Now  $x_{t+1} = x_t - [g''(x_t)]^{-1} g'(x_t)$

$$(x_{t+1} - x^*) = (x_t - x^*)^2 \frac{g'''(x^b)}{2g''(x_t)}$$

$$\text{So } \mathcal{E}_{t+1} = \mathcal{E}_t^2 \frac{g'''(x^b)}{2g''(x_t)}$$

Now,  $\left| \frac{g'''(x^b)}{2g''(x_t)} \right|$  is continuous for  $(x^b, x_t) \in \underbrace{(x^* \pm \varepsilon) \times (x^* \pm \varepsilon)}_{\text{Cartesian product}}$

So continuous on compact  $S = [x^* \pm \delta] \times [x^* \pm \delta]$   
if  $\delta < \varepsilon \implies$  has a maximum on  $S$ , say  $M$ .

$$\therefore |\mathcal{E}_{t+1}| \leq \mathcal{E}_t^2 M \quad \text{if } |\mathcal{E}_t| < \delta$$

and  $|\varepsilon_{t+1}| \leq \frac{1}{2} |\varepsilon_t|$  if  $|\varepsilon_t| < \min\{\delta, \frac{1}{2M}\}$

Therefore  $|\varepsilon_t| \rightarrow 0$  as  $t \rightarrow \infty$

and  $x_t \rightarrow x^*$

if we start within  $\min\{\delta, \frac{1}{2M}\}$  of  $x^*$ .

---

---

When can we be sure that Newton's method will work?

NMC 1:

Suppose: 1)  $x_*$  <sup>such that</sup>  $\Rightarrow g'(x_*) = 0$

2)  $g' \in C^2$  near  $x_*$

3)  $g''(x_*) \neq 0$

$x_*$  is a root of  $g'$

$C^2$ : twice continuously differentiable

Then  $\exists \epsilon > 0$

not the same as the one we started with above

my notation for an interval

$\exists$  starting from any  $x_0 \in (x_* \pm \epsilon)$

such that

Newton's Method will converge to  $x_*$

Proof: Above.



When can we be sure that Newton's method will work?

NMC 1:

Suppose: 1)  $x_* \rightarrow g'(x_*) = 0$

such that

$x_*$  is a root of  $g'$

2)  $g' \in C^2$  near  $x_*$

$C^2$ : twice continuously differentiable

there exists

3)  $g''(x_*) \neq 0$

Then  $\exists \epsilon > 0$

my notation for an interval

$\exists$  starting from any  $x_0 \in (x_* \pm \epsilon)$

such that

Newton's Method will converge to  $x_*$

Proof: Above.

NM C2: If  $g \in C^3$ ,  $g''(x_*) < 0$   
and  $g$  concave down  
then NM converges from anywhere.

---

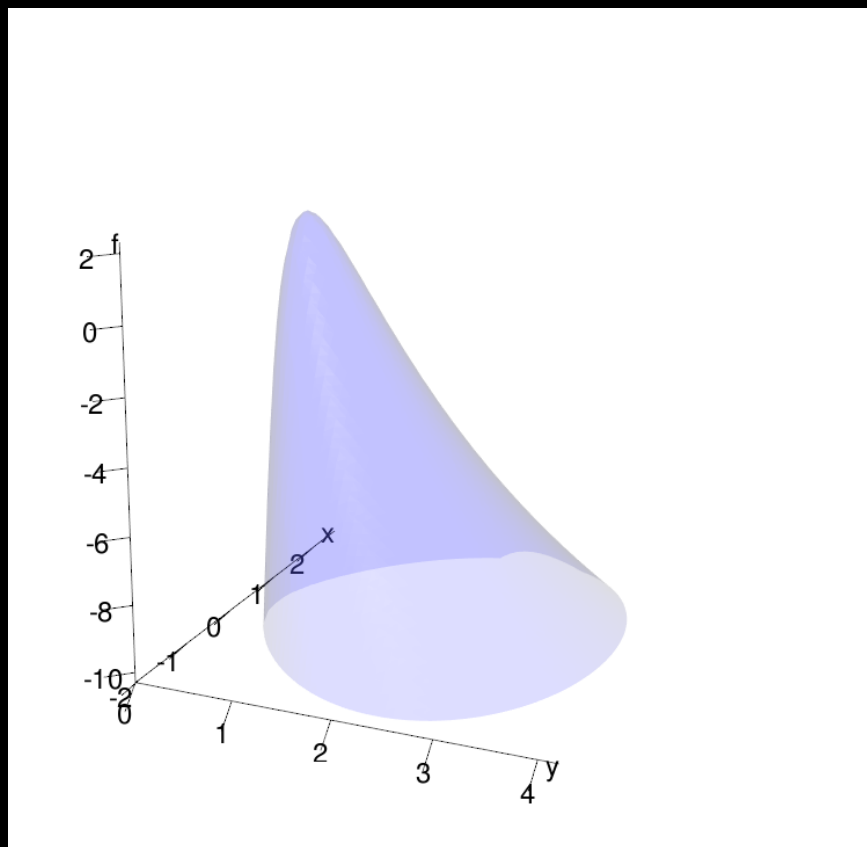
So alternative strategy:

Reparametrize  $g$  to make it concave down

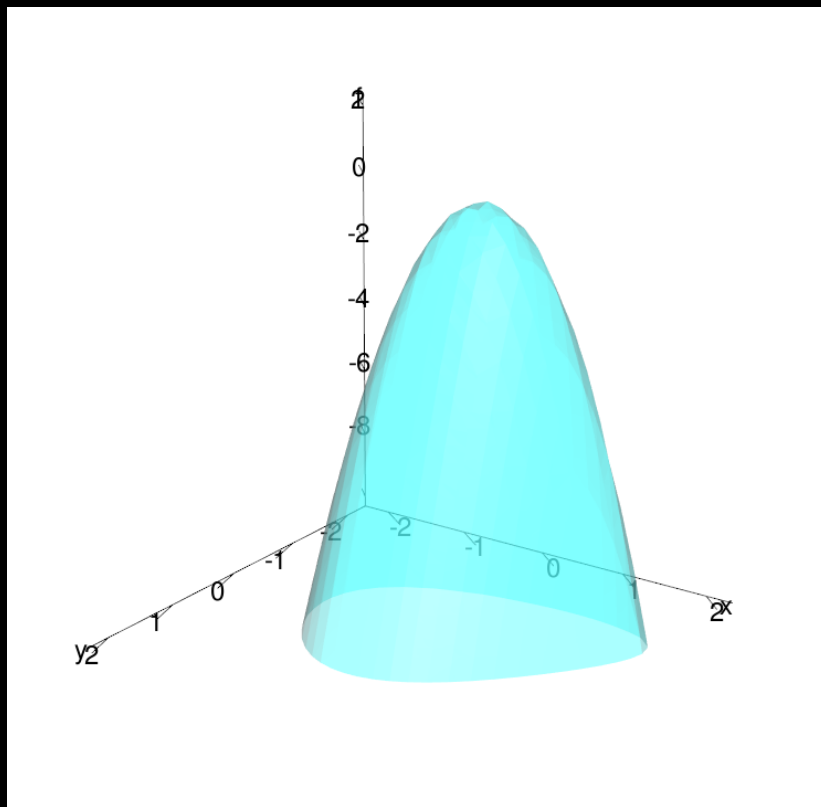
e.g. Instead of  $f(x | \mu, \sigma)$

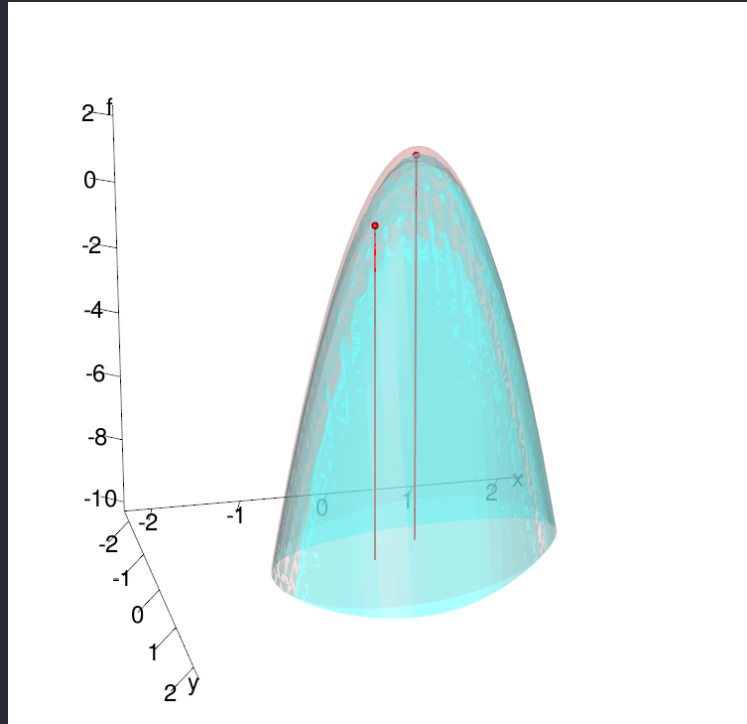
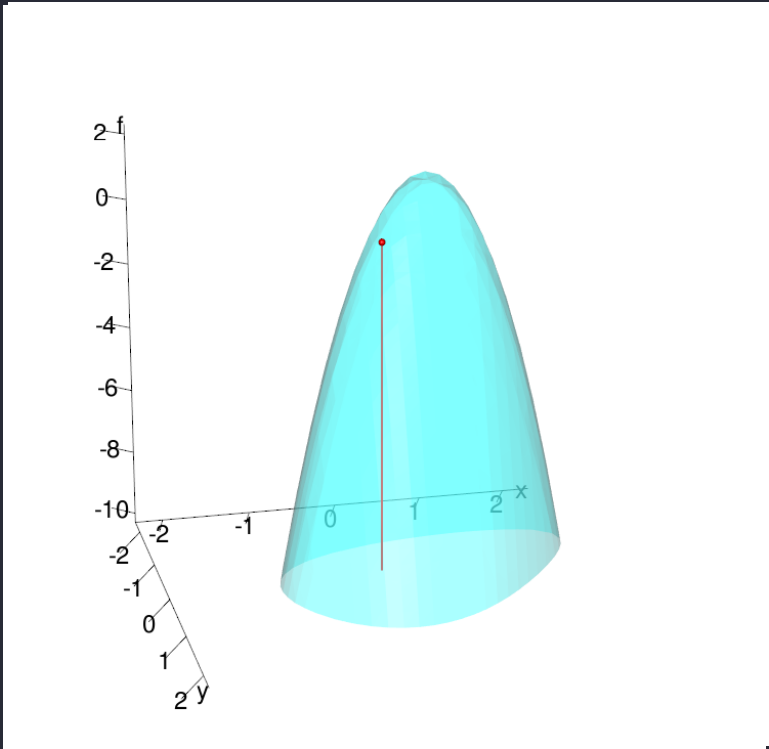
use  $f(x | \theta, \psi)$  where  $\theta = \mu/\sigma$   $\psi = \log \sigma$

Using  $\mu, \sigma$ :

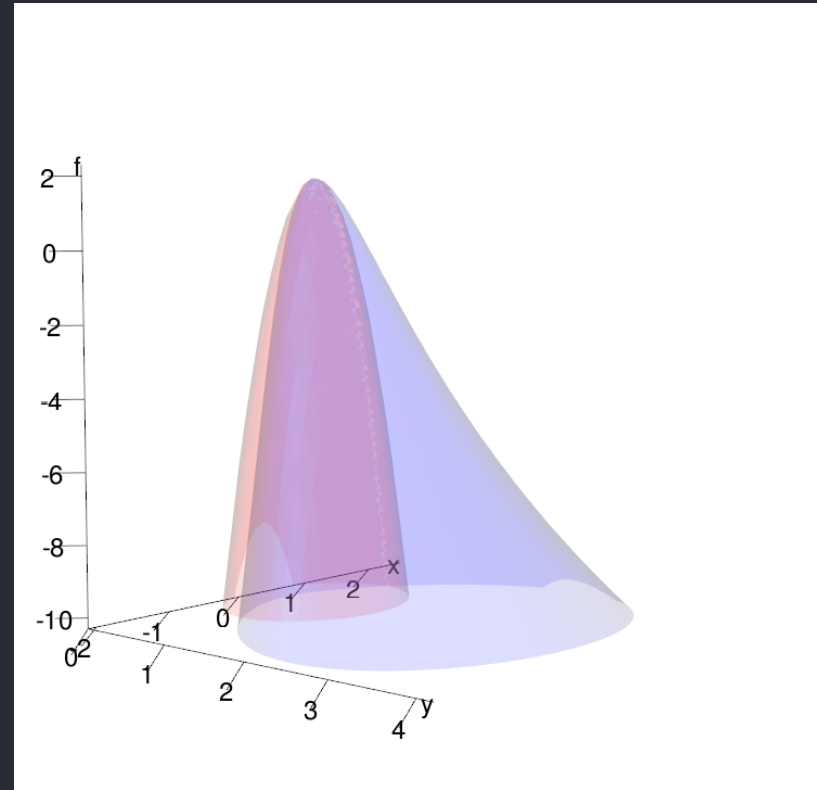
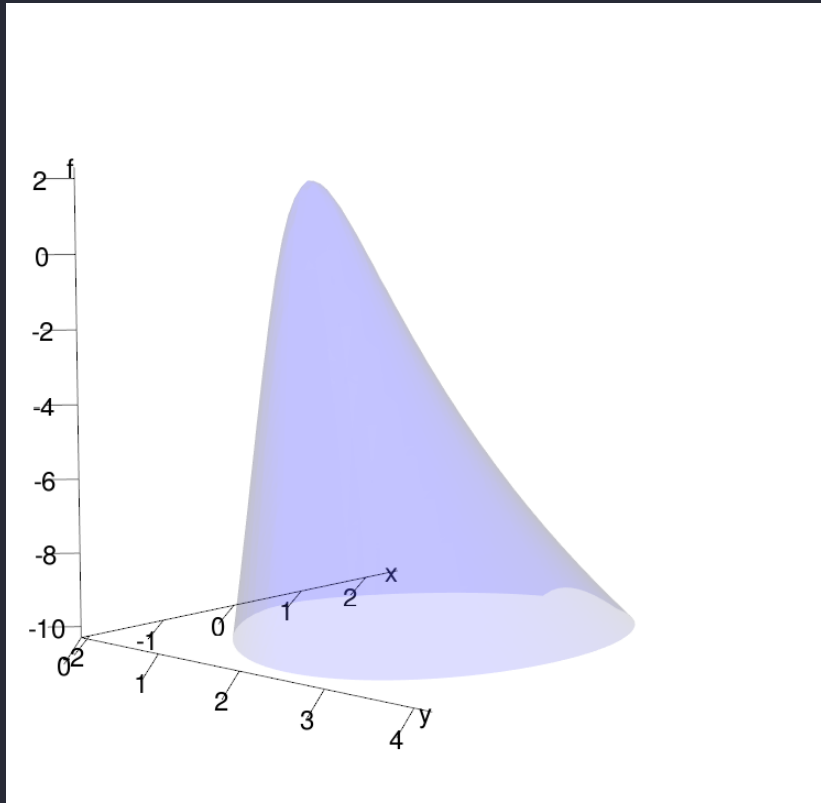


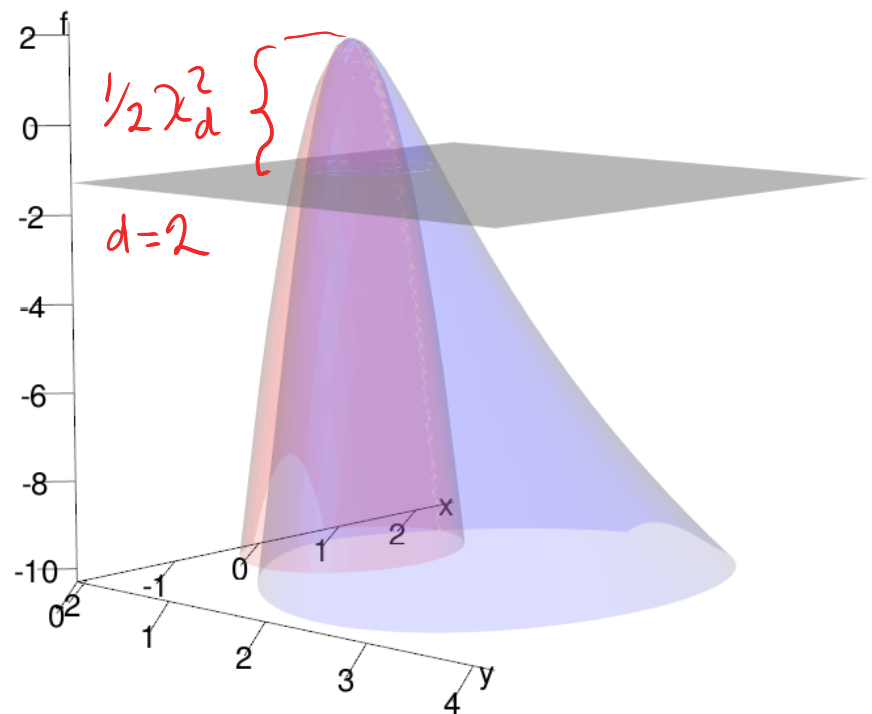
Using  $\theta, \psi$ :





# LRT vs. Wald Confidence regions





LR 95% CR

Wald 95% CR

- LR CR is equivariant under reparametrization.
- Wald CR is NOT but is much closer to LR CR if likelihood is reparametrized to be more quadratic.

