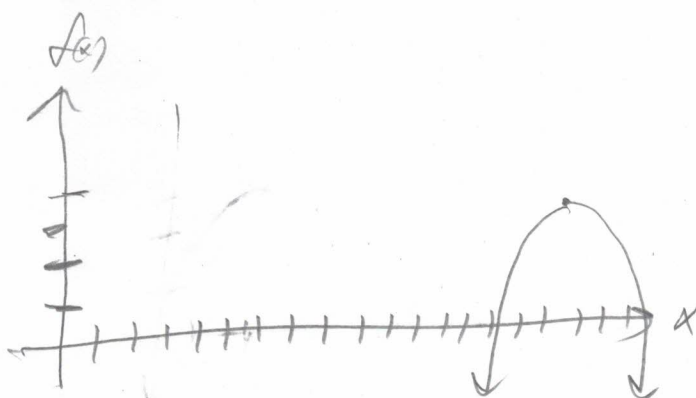


MATH 390/690 LEC 23. (FINAL!!)

Prob: let $f(x) = 3 - (x-17)^2$



$$f_* = \max_{x \in \mathbb{R}} \{f(x)\} = 3$$

$$x_* = \operatorname{argmax}_{x \in \mathbb{R}} \{f(x)\} := \{x: f(x) = f_*\} = \{x: f(x) = 3\} = \{17\}$$

How to find x_* ? Set $f'(x) = 0$ and solve. Include if $f''(x) < 0$.

$$f'(x) = -2(x-17) = 0 \Rightarrow x-17 = 0 \Rightarrow x = 17, f''(x) = -2 \checkmark$$

Note argmax is affected by strictly increasing function of x . E.g. if

$$\text{e.g. } g(x) = q f(x) \quad \text{st } q > 0 \Rightarrow \operatorname{argmax} \{f(x)\} = \operatorname{argmax} \{g(x)\}$$

$$\text{e.g. Also } g(x) = \ln(f(x))$$

If $f(x)$ is a PDF, or $p(x)$ is a PMF for rv X

$$\text{Mode}(X) \stackrel{\text{com}}{:=} \operatorname{argmax}_{x \in S_X} \{f(x)\} =$$

$$\stackrel{\text{dim}}{=} \operatorname{argmax}_{x \in S_X} \{p(x)\} = \operatorname{argmax}_{x \in S_X} \{f(x)\} = \operatorname{argmax}_{x \in S_X} \{\ln(f(x))\}$$

Finding modes is important in estimation schemes (MATH 391)

Also, it is a general useful problem called "optimization". Big field!

e.g. $X \sim N(\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$

$$\propto e^{-\frac{1}{2\sigma^2}x^2 + \frac{x\mu}{\sigma^2} - \frac{\mu^2}{2\sigma^2}}$$

$$\propto e^{-\frac{1}{2\sigma^2}x^2 + \frac{\mu}{\sigma^2}x} = h(x) \Rightarrow \ln(h(x)) = \frac{\mu}{\sigma^2}x - \frac{1}{2\sigma^2}x^2$$

$$\text{Mode}(X) = \underset{x}{\text{argmax}} \left\{ \frac{\mu}{\sigma^2}x - \frac{1}{2\sigma^2}x^2 \right\}$$

Since X is continuous...

$$\Rightarrow g'(x) = \frac{\mu}{\sigma^2} - \frac{1}{\sigma^2}x \stackrel{\text{set}}{=} 0 \Rightarrow x_0 = \mu$$

$g(x)$

Note: we skipped Poisson process, Skellam distribution

end of course

Let \vec{x} be a vector rv with arbitrary $\vec{\mu}$, Σ matrix

$$\text{Mode}(\vec{x}) = \underset{\vec{x} \in \mathcal{X}}{\text{argmax}} \{ \ln(h_{\vec{x}}(\vec{x})) \} \quad \text{HARDER!}$$

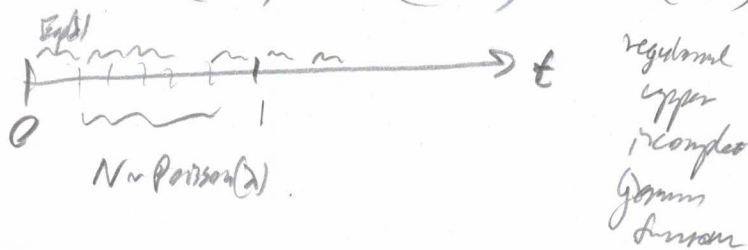
MATH 202 gives you tools to answer this under some conditions but in the real world, those conditions don't hold.

In fact, usually it is impossible to answer, so instead it is approximated as best as possible.

Especially since $\dim(\vec{x})$ is large. In modern statistical models this can be in the 1000's. Chat GPT has $\dim(\vec{x}) = 100$ trillion!

$$T_k \sim \text{Erlang}(k, \lambda), N \sim \text{Poisson}(\lambda)$$

$$P(T_k > 1) = P(N \leq k-1) = Q(k, \lambda)$$



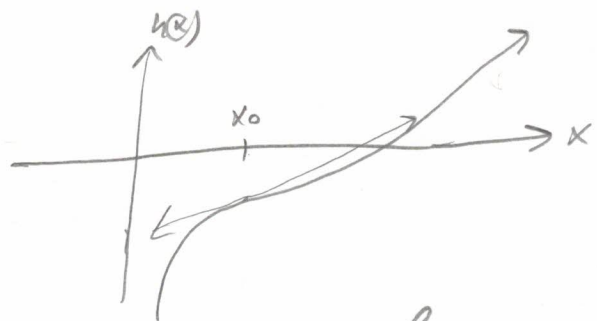
the Newton-Raphson alg.

3

Let's start with an approximation algorithm in one dimension

Let $g(x) = \ln(k(x))$ which is assumed differentiable

Let $h(x) := g'(x)$. We want to solve for $h(x) = 0$.



Step 0: Define ϵ , a tolerance e.g. $\epsilon = .01$ (where the solution will not differ from the approx sol. by ϵ .)

Step 1: guess the solution at $x = x_0$

Step 2: Draw tangent line at $h(x_0)$, call it $l(x_0)$

Solve for x where $l(x_0) = 0$ (i.e. find its x intercept)

And call its solution x_1 .
$$y - y_0 = m(x - x_0) \Rightarrow y = mx + mx_0 + y_0 \stackrel{\text{set } 0}{\Rightarrow} mx = -mx_0 - y_0 \Rightarrow x_1 = x_0 - \frac{y_0}{m} = x_0 - \frac{h(x_0)}{h'(x_0)}$$

Step 3: Repeat Step 2, until $|x_i - x_{i-1}| < \epsilon$ "stopping rule"

E.g. $X \sim \text{Beta}(\overset{\alpha}{17}, \overset{\beta}{37})$. $\text{Mode}(X) = \frac{\alpha - 1}{\alpha + \beta - 2} = \frac{17 - 1}{17 + 37 - 2} = .30769$

Try Newton Raphson. $f(x) = \frac{1}{B(17,37)} x^{16} (1-x)^{36-1} \Rightarrow k(x) = x^{16} (1-x)^{36}$
 $\Rightarrow g(x) = \ln(k(x)) = 16 \ln(x) + 36 \ln(1-x)$
 $\Rightarrow h(x) = \frac{16}{x} - \frac{36}{1-x}$
 $\Rightarrow h'(x) = -\frac{16}{x^2} - \frac{36}{(1-x)^2}$

Step 1: good guess: $x_0 = 0.5$

Step 2: $h'(x_0) = -\frac{16}{0.5^2} - \frac{36}{(0.5)^2} = -208$
 $h(x_0) = \frac{16}{0.5} - \frac{36}{0.5} = -10$
 $\Rightarrow x_1 = 0.5 - \frac{-10}{-208} = 0.9519$ $|x_1 - x_0| = .099 \not< \epsilon$
continue

Step 2: $h'(x_1) = -\frac{16}{.9519^2} - \frac{36}{(1-.9519)^2} = -30.215$
 $h(x_1) = \frac{16}{.9519} - \frac{36}{(1-.9519)} = -198.10$
 $\Rightarrow x_2 = .9519 - \frac{-198.10}{-30.215} = .2991$
 $|x_2 - x_1| = .152 \not< \epsilon$
continue

Step 23: $h'(x_2) = -252.1$
 $h(x_2) = 2.122$

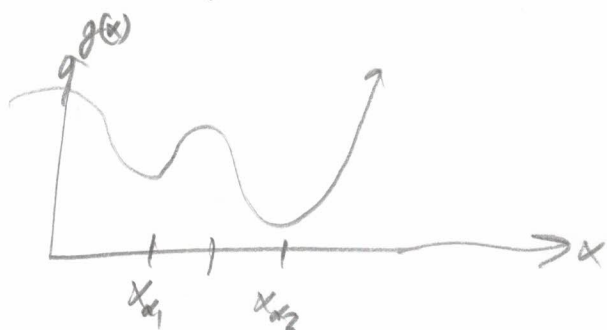
$x_3 = .2991 - \frac{2.122}{-252.1} = .30755$

$(x_3 - x_2) = .008 < \epsilon \Rightarrow \underline{\underline{\text{STOP}}}$

Very close to true answer!!

Advantage: we never overshoot the minimum by stepping "too far".

Major pitfall: if function has multiple minimums: x_{k1}, x_{k2}
 one being lower than the other, $g(x_{k2}) < g(x_{k1})$, different



$x_0 \Rightarrow$ Converges to x_{k1} | $x_0 \Rightarrow$ Converges to x_{k2}

Value of $x_0 \Rightarrow$
 different answers.

Other problems: slow! It
 Can get faster if you use
 2nd derivative of h as well
 Also, sometimes you need to examine the
 function itself heuristically \Rightarrow HARD

Can this be extended to dimensions ⁿ greater than 1? Yes..

let $g(\vec{x}) := h_n(\vec{x}) \in \mathbb{R}^n$

Now, the "derivative" of g is derivative of each dimension

called the "gradient"

$$h(\vec{x}) := \vec{\nabla} g(\vec{x}) := \begin{bmatrix} \frac{\partial g}{\partial x_1} |_{\vec{x}} \\ \vdots \\ \frac{\partial g}{\partial x_n} |_{\vec{x}} \end{bmatrix} \in \mathbb{R}^n$$

What is the analogue of the $h'(\vec{x})$ here? This is called the

Hessian

$$H(\vec{x}) := \begin{bmatrix} \frac{\partial^2 g}{\partial x_1^2} |_{\vec{x}} & \frac{\partial^2 g}{\partial x_1 \partial x_2} |_{\vec{x}} & \dots & \frac{\partial^2 g}{\partial x_1 \partial x_n} |_{\vec{x}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 g}{\partial x_n \partial x_1} |_{\vec{x}} & \frac{\partial^2 g}{\partial x_n \partial x_2} |_{\vec{x}} & \dots & \frac{\partial^2 g}{\partial x_n^2} |_{\vec{x}} \end{bmatrix}$$

Step 0: ϵ decided step 1: x_0 decided

Step 2: Now the update step becomes

$$\vec{x}_1 = \vec{x}_0 - H(\vec{x}_0)^{-1} \vec{\nabla} h(\vec{x}_0)$$

Step 3: Continue until $|\vec{x}_i - \vec{x}_{i-1}| < \epsilon$.

All the problems in the 1-d case apply here
and are a lot worse! Additionally: the Hessian is expensive to
compute and solution is extracted from data and can be biased!

$g(\vec{x}) := -\ln(K(\vec{x}))$, so how to find argmin $(\ln(K(\vec{x})))$ is

find argmin $\{g(\vec{x})\}$

Is there a less expensive algorithm? How about not using Hessian?

Defining ϵ as step magnitude and using as update the:

$$\vec{x}_i = \vec{x}_{i-1} - \epsilon \vec{\nabla} h(\vec{x}_{i-1})$$

AKA "learning rate"

descent is the direction of the gradient \Rightarrow gradient descent

If ϵ is constant throughout the iterations, we
can move too slow or we can move too fast

Stochastic Gradient Descent: When we use data to
compute $\vec{\nabla} h(\vec{x})$, we don't use all the data, we take a
random subset \Rightarrow better performance and better speed