

# PDE-Constrained and Nonsmooth Optimization

Frank E. Curtis

COR@L Seminar

October 1, 2009

# Outline

## PDE-Constrained Optimization

- Introduction

- Newton's method

- Inexactness

- Results

- Summary and future work

## Nonsmooth Optimization

- Sequential quadratic programming (SQP)

- Gradient sampling (GS)

- SQP-GS

- Results

- Summary and future work

## Conclusion

## Introduction

### Model 1

6

### Sequential quadratic

Gradient coupling (CG)

SOP, CS

D. K.

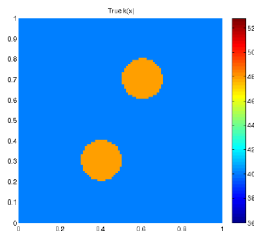
6

## PDE-constrained optimization

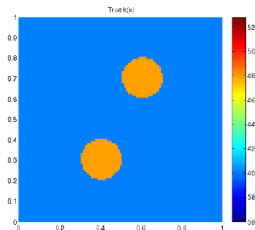
$$\begin{array}{ll} \min & f(x) \\ \text{s.t.} & c_{\mathcal{E}}(x) = 0 \\ & c_{\mathcal{I}}(x) \geq 0 \end{array}$$

$$\begin{array}{ll} \min & f(x) \\ \text{s.t.} & c_{\mathcal{E}}(x) = 0 \quad (PDE) \\ & c_{\mathcal{I}}(x) \geq 0 \end{array}$$

Problem is infinite-dimensional

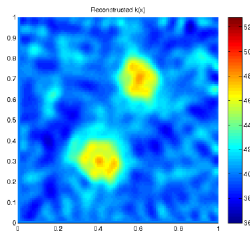
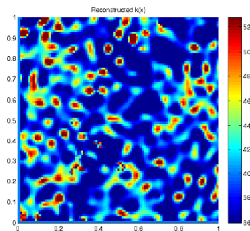


1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100



$$\min_{y,k} \frac{1}{2} \sum_j \sum_m (y_j(x_m) - y_{j,m})^2 + \alpha(\beta \|k\|_{L^2(\Omega)}^2 + \|\nabla k\|_{(L^2(\Omega))^n}^2)$$

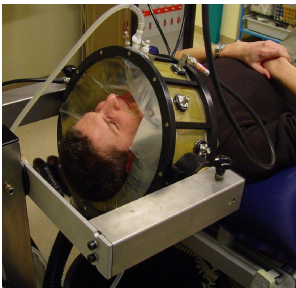
$$\text{s.t.} \begin{cases} \Delta y_j(x) + S(x)k(x)^2 y_j(x) - S(x)(k_0^2 - k(x)^2) y_j^i = 0, & x \text{ in } \Omega \\ y_j = 0, & x \text{ on } \partial\Omega \\ l(x) \leq k(x) \leq u(x), & x \text{ in } \Omega \end{cases}$$

[illegible]

$$\min_{y,k} \frac{1}{2} \sum_j \sum_m (y_j(x_m) - y_{j,m})^2 + \alpha (\beta \|k\|_{L^2(\Omega)}^2 + \|\nabla k\|_{(L^2(\Omega))^n}^2)$$

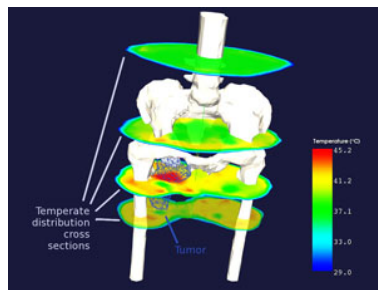
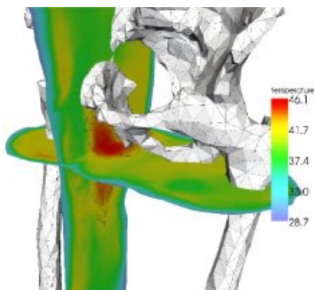
$$\text{s.t.} \quad \begin{cases} \Delta y_j(x) + S(x)k(x)^2 y_j(x) - S(x)(k_0^2 - k(x)^2) y_j^i = 0, & x \text{ in } \Omega \\ y_j = 0, & x \text{ on } \partial\Omega \\ l(x) \leq k(x) \leq u(x), & x \text{ in } \Omega \end{cases}$$





1. **Introduction**

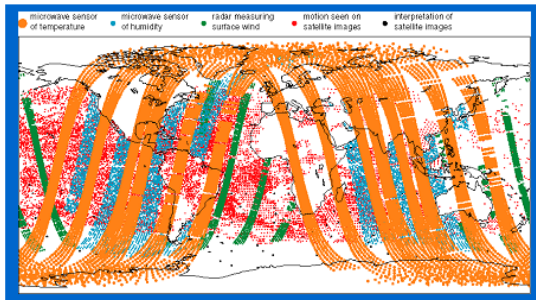
- ▶ Computer modeling can be used to help plan the therapy for each patient, and it opens the door for numerical optimization
- ▶ The goal is to heat the tumor to a target temperature of  $43^{\circ}\text{C}$  while minimizing damage to nearby cells





- ▶ If the initial state of the atmosphere (temperatures, pressures, wind patterns, humidities) were known at a certain point in time, then an accurate forecast could be obtained by integrating atmospheric model equations forward in time
- ▶ Flow described by Navier-Stokes and further sophistications of atmospheric physics and dynamics

- ▶ Each observation is subject to error
- ▶ Nonuniformly distributed around the globe (satellite paths, densely-populated areas)



### Model 1

6

### Sequential quadratic

### Gradient sampling (GS)

SOD, CF

D. 1.4.

6

# Nonlinear equations

- ▶ Newton's method

$$\boxed{\mathcal{F}(x) = 0} \Rightarrow \boxed{\nabla \mathcal{F}(x_k) d_k = -\mathcal{F}(x_k)}$$

- ▶ Judge progress by the merit function

$$\phi(x) \triangleq \frac{1}{2} \|\mathcal{F}(x)\|^2$$

- ▶ Direction is one of descent since

$$\nabla \phi(x_k)^T d_k = \mathcal{F}(x_k)^T \nabla \mathcal{F}(x_k) d_k = -\|\mathcal{F}(x_k)\|^2 < 0$$

(Note the **consistency** between the step computation and merit function!)

- $$\mathcal{L}(x, \lambda) \triangleq f(x) + \lambda^T c(x)$$

- ▶ Simply minimizing

$$\varphi(x, \lambda) = \frac{1}{2} \|\mathcal{F}(x, \lambda)\|^2 = \frac{1}{2} \left\| \begin{bmatrix} \nabla f(x) + \nabla c(x)\lambda \\ c(x) \end{bmatrix} \right\|^2$$

is generally inappropriate for constrained optimization

- ▶ We use the merit function

$$\phi(x; \pi) \triangleq f(x) + \pi \|c(x)\|$$

where  $\pi$  is a penalty parameter



6.  $\frac{1}{2} \times \frac{1}{2} = \frac{1}{4}$

$$\min (x-1)^2, \text{ s.t. } x=0 \quad \text{i.e.} \quad \phi(x; \pi) = (x-1)^2 + \pi|x|$$

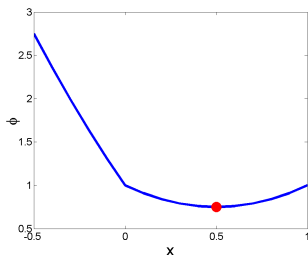


Figure:  $\pi = 1$

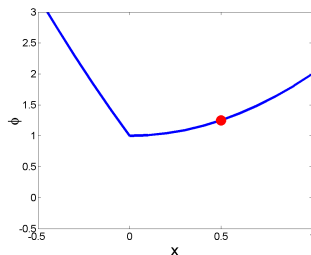


Figure:  $\pi = 2$

for  $k = 0, 1, 2$ 

$$\begin{bmatrix} H(x_k, \lambda_k) & \nabla c(x_k) \\ \nabla c(x_k)^T & 0 \end{bmatrix} \begin{bmatrix} d_k \\ \delta_k \end{bmatrix} = - \begin{bmatrix} \nabla f(x_k) + \nabla c(x_k) \lambda_k \\ c(x_k) \end{bmatrix}$$

- 1997, 1998, 1999, 2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018, 2019, 2020, 2021, 2022, 2023, 2024, 2025, 2026, 2027, 2028, 2029, 2030, 2031, 2032, 2033, 2034, 2035, 2036, 2037, 2038, 2039, 2040, 2041, 2042, 2043, 2044, 2045, 2046, 2047, 2048, 2049, 2050, 2051, 2052, 2053, 2054, 2055, 2056, 2057, 2058, 2059, 2060, 2061, 2062, 2063, 2064, 2065, 2066, 2067, 2068, 2069, 2070, 2071, 2072, 2073, 2074, 2075, 2076, 2077, 2078, 2079, 2080, 2081, 2082, 2083, 2084, 2085, 2086, 2087, 2088, 2089, 2090, 2091, 2092, 2093, 2094, 2095, 2096, 2097, 2098, 2099, 2100, 2101, 2102, 2103, 2104, 2105, 2106, 2107, 2108, 2109, 2110, 2111, 2112, 2113, 2114, 2115, 2116, 2117, 2118, 2119, 2120, 2121, 2122, 2123, 2124, 2125, 2126, 2127, 2128, 2129, 2130, 2131, 2132, 2133, 2134, 2135, 2136, 2137, 2138, 2139, 2140, 2141, 2142, 2143, 2144, 2145, 2146, 2147, 2148, 2149, 2150, 2151, 2152, 2153, 2154, 2155, 2156, 2157, 2158, 2159, 2160, 2161, 2162, 2163, 2164, 2165, 2166, 2167, 2168, 2169, 2170, 2171, 2172, 2173, 2174, 2175, 2176, 2177, 2178, 2179, 2180, 2181, 2182, 2183, 2184, 2185, 2186, 2187, 2188, 2189, 2190, 2191, 2192, 2193, 2194, 2195, 2196, 2197, 2198, 2199, 2200, 2201, 2202, 2203, 2204, 2205, 2206, 2207, 2208, 2209, 2210, 2211, 2212, 2213, 2214, 2215, 2216, 2217, 2218, 2219, 2220, 2221, 2222, 2223, 2224, 2225, 2226, 2227, 2228, 2229, 2230, 2231, 2232, 2233, 2234, 2235, 2236, 2237, 2238, 2239, 2240, 2241, 2242, 2243, 2244, 2245, 2246, 2247, 2248, 2249, 2250, 2251, 2252, 2253, 2254, 2255, 2256, 2257, 2258, 2259, 2260, 2261, 2262, 2263, 2264, 2265, 2266, 2267, 2268, 2269, 2270, 2271, 2272, 2273, 2274, 2275, 2276, 2277, 2278, 2279, 2280, 2281, 2282, 2283, 2284, 2285, 2286, 2287, 2288, 2289, 2290, 2291, 2292, 2293, 2294, 2295, 2296, 2297, 2298, 2299, 2300, 2301, 2302, 2303, 2304, 2305, 2306, 2307, 2308, 2309, 2310, 2311, 2312, 2313, 2314, 2315, 2316, 2317, 2318, 2319, 2320, 2321, 2322, 2323, 2324, 2325, 2326, 2327, 2328, 2329, 2330, 2331, 2332, 2333, 2334, 2335, 2336, 2337, 2338, 2339, 2340, 2341, 2342, 2343, 2344, 2345, 2346, 2347, 2348, 2349, 2350, 2351, 2352, 2353, 2354, 2355, 2356, 2357, 2358, 2359, 2360, 2361, 2362, 2363, 2364, 2365, 2366, 2367, 2368, 2369, 2370, 2371, 2372, 2373, 2374, 2375, 2376, 2377, 2378, 2379, 2380, 2381, 2382, 2383, 2384, 2385, 2386, 2387, 2388, 2389, 2390, 2391, 2392, 2393, 2394, 2395, 2396, 2397, 2398, 2399, 2400, 2401, 2402, 2403, 2404, 2405, 2406, 2407, 2408, 2409, 2410, 2411, 2412, 2413, 2414, 2415, 2416, 2417, 2418, 2419, 2420, 2421, 2422, 2423, 2424, 2425, 2426, 2427, 2428, 2429, 2430, 2431, 2432, 2433, 2434, 2435, 2436, 2437, 2438, 2439, 2440, 2441, 2442, 2443, 2444, 2445, 2446, 2447, 2448, 2449, 2450, 2451, 2452, 2453, 2454, 2455, 2456, 2457, 2458, 2459, 2460, 2461, 2462, 2463, 2464, 2465, 2466, 2467, 2468, 2469, 2470, 2471, 2472, 2473, 2474, 2475, 2476, 2477, 2478, 2479, 2480, 2481, 2482, 2483, 2484, 2485, 2486, 2487, 2488, 2489, 2490, 2491, 2492, 2493, 2494, 2495, 2496, 2497, 2498, 2499, 2500, 2501, 2502, 2503, 2504, 2505, 2506, 2507, 2508, 2509, 2510, 2511, 2512, 2513, 2514, 2515, 2516, 2517, 2518, 2519, 2520, 2521, 2522, 2523, 2524, 2525, 2526, 2527, 2528, 2529, 2530, 2531, 2532, 2533, 2534, 2535, 2536, 2537, 2538, 2539, 2540, 2541, 2542, 2543, 2544, 2545, 2546, 2547, 2548, 2549, 2550, 2551, 2552, 2553, 2554, 2555, 2556, 2557, 2558, 2559, 2560, 2561, 2562, 2563, 2564, 2565, 2566, 2567, 2568, 2569, 2570, 2571, 2572, 2573, 2574, 2575, 2576, 2577, 2578, 2579, 2580, 2581, 2582, 2583, 2584, 2585, 2586, 2587, 2588, 2589, 2590, 2591, 2592, 2593, 2594, 2595, 2596, 2597, 2598, 2599, 2600, 2601, 2602, 2603, 2604, 2605, 2606, 2607, 2608, 2609, 2610, 2611, 2612, 2613, 2614, 2615, 2616, 2617, 2618, 2619, 2620, 2621, 2622, 2623, 2624, 2625, 2626, 2627, 2628, 2629, 2630, 2631, 2632, 2633, 2634, 2635, 2636, 2637, 2638, 2639, 2640, 2641, 2642, 2643, 2644, 2645, 2646, 2647, 2648, 2649, 2650, 2651, 2652, 2653, 2654, 2655, 2656, 2657, 2658, 2659, 2660, 2661, 2662, 2663, 2664, 2665, 2666, 2667, 2668, 2669, 2670, 2671, 2672, 2673, 2674, 2675, 2676, 2677, 2678, 26

$$\phi(\mathbf{x}_k + \alpha_k \mathbf{d}_k; \pi_k) \leq \phi(\mathbf{x}_k; \pi_k) + \eta \alpha_k D\phi_k(\mathbf{d}_k; \pi_k)$$

- **Update** iterate  $(x_{k+1}, \lambda_{k+1}) \leftarrow (x_k, \lambda_k) + \alpha_k(d_k, \delta_k)$

100

71

- ▶ **(Regularity)**  $\nabla c(x_k)^T$  has full row rank with singular values bounded below by a positive constant
- ▶ **(Convexity)**  $u^T H(x_k, \lambda_k) u \geq \mu \|u\|^2$  for  $\mu > 0$  for all  $u \in \mathbb{R}^n$  satisfying  $u \neq 0$  and  $\nabla c(x_k)^T u = 0$

## (11) (1077)

$$\lim_{k \rightarrow \infty} \left\| \begin{bmatrix} \nabla f(x_k) + \nabla c(x_k) \lambda_k \\ c(x_k) \end{bmatrix} \right\| = 0$$

# Outline

## PDE-Constrained Optimization

Introduction

Newton's method

**Inexactness**

Results

Summary and future work

## Nonsmooth Optimization

Sequential quadratic programming (SQP)

Gradient sampling (GS)

SQP-GS

Results

Summary and future work

## Conclusion

- ▶ Computational issues:
  - ▶ Large matrices to be stored
  - ▶ Large matrices to be factored
- ▶ Algorithmic issues:
  - ▶ The problem may be nonconvex
  - ▶ The problem may be ill-conditioned
- ▶ Computational/Algorithmic issues:
  - ▶ No matrix factorizations makes difficulties more difficult

- $$\nabla \mathcal{F}(x_k) d_k = -\mathcal{F}(x_k) + r_k$$
- requiring (Dembo, Eisenstat, Steihaug (1982))

- $$\phi(x) \triangleq \frac{1}{2} \|\mathcal{F}(x)\|^2$$

- $$\phi(x_k)^T d_k = \mathcal{F}(x_k)^T \nabla \mathcal{F}(x_k) d_k = -\|\mathcal{F}(x_k)\|^2 + \mathcal{F}(x_k)^T r_k \leq (\kappa-1)\|\mathcal{F}(x_k)\|^2 < 0$$

\_\_\_\_\_





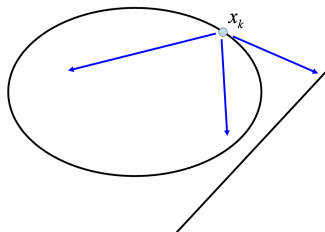
# Termination test 1

The search direction  $(d_k, \delta_k)$  is **acceptable** if

$$\left\| \begin{bmatrix} \rho_k \\ r_k \end{bmatrix} \right\| \leq \kappa \left\| \begin{bmatrix} \nabla f(x_k) + \nabla c(x_k) \lambda_k \\ c(x_k) \end{bmatrix} \right\|, \quad \kappa \in (0, 1)$$

and if for  $\pi_k = \pi_{k-1}$  and some  $\sigma \in (0, 1)$  we have

$$\Delta m(d_k; \pi_k) \geq \underbrace{\max\left\{\frac{1}{2}d_k^T H(x_k, \lambda_k)d_k, 0\right\} + \sigma\pi_k \max\{\|c(x_k)\|, \|r_k\| - \|c(x_k)\|\}}_{\geq 0 \text{ for any } d}$$

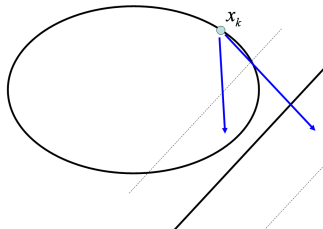


## Termination test 2

The search direction  $(d_k, \delta_k)$  is **acceptable** if

$$\|\rho_k\| \leq \beta \|c(x_k)\|, \quad \beta > 0$$

$$\text{and} \quad \|r_k\| \leq \epsilon \|c(x_k)\|, \quad \epsilon \in (0, 1)$$



Increasing the penalty parameter  $\pi$  then yields

$$\Delta m(d_k; \pi_k) \geq \underbrace{\max\left\{\frac{1}{2} d_k^T H(x_k, \lambda_k) d_k, 0\right\} + \sigma \pi_k \|c(x_k)\|}_{\geq 0 \text{ for any } d}$$

# Algorithm 1: Inexact Newton for optimization

(Byrd, Curtis, Nocedal (2008))

for  $k = 0, 1, 2, \dots$

- Iteratively solve

$$\begin{bmatrix} H(x_k, \lambda_k) & \nabla c(x_k) \\ \nabla c(x_k)^T & 0 \end{bmatrix} \begin{bmatrix} d_k \\ \delta_k \end{bmatrix} = - \begin{bmatrix} \nabla f(x_k) + \nabla c(x_k) \lambda_k \\ c(x_k) \end{bmatrix}$$

until termination test 1 or 2 is satisfied

- If only termination test 2 is satisfied, increase  $\pi$  so

$$\pi_k \geq \max \left\{ \pi_{k-1}, \frac{\nabla f(x_k)^T d_k + \max\{\frac{1}{2} d_k^T H(x_k, \lambda_k) d_k, 0\}}{(1 - \tau)(\|c(x_k)\| - \|r_k\|)} \right\}$$

- Backtrack from  $\alpha_k \leftarrow 1$  to satisfy

$$\phi(x_k + \alpha_k d_k; \pi_k) \leq \phi(x_k; \pi_k) - \eta \alpha_k \Delta m(d_k; \pi_k)$$

- Update iterate  $(x_{k+1}, \lambda_{k+1}) \leftarrow (x_k, \lambda_k) + \alpha_k (d_k, \delta_k)$

# Convergence of Algorithm 1

## Assumption

The sequence  $\{(x_k, \lambda_k)\}$  is contained in a convex set  $\Omega$  over which  $f$ ,  $c$ , and their first derivatives are bounded and Lipschitz continuous. Also,

- ▶ (**Regularity**)  $\nabla c(x_k)^T$  has full row rank with singular values bounded below by a positive constant
- ▶ (**Convexity**)  $u^T H(x_k, \lambda_k) u \geq \mu \|u\|^2$  for  $\mu > 0$  for all  $u \in \mathbb{R}^n$  satisfying  $u \neq 0$  and  $\nabla c(x_k)^T u = 0$

## Theorem

(Byrd, Curtis, Nocedal (2008)) The sequence  $\{(x_k, \lambda_k)\}$  yields the limit

$$\lim_{k \rightarrow \infty} \left\| \begin{bmatrix} \nabla f(x_k) + \nabla c(x_k) \lambda_k \\ c(x_k) \end{bmatrix} \right\| = 0$$

# Handling nonconvexity and rank deficiency

- ▶ There are two assumptions we aim to drop:
  - ▶ (*Regularity*)  $\nabla c(x_k)^T$  has full row rank with singular values bounded below by a positive constant
  - ▶ (*Convexity*)  $u^T H(x_k, \lambda_k) u \geq \mu \|u\|^2$  for  $\mu > 0$  for all  $u \in \mathbb{R}^n$  satisfying  $u \neq 0$  and  $\nabla c(x_k)^T u = 0$

e.g., the problem is not regular if it is **infeasible**, and it is not convex if there are **maximizers and/or saddle points**

- ▶ Without them, Algorithm 1 may stall or may not be well-defined

# No factorizations means no clue

- ▶ We might not **store** or **factor**

$$\begin{bmatrix} H(x_k, \lambda_k) & \nabla c(x_k) \\ \nabla c(x_k)^T & 0 \end{bmatrix}$$

so we might not know if the problem is **nonconvex** or **ill-conditioned**

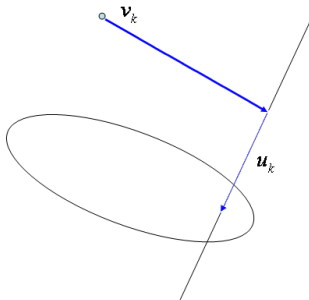
- ▶ Common practice is to perturb the matrix to be

$$\begin{bmatrix} H(x_k, \lambda_k) + \xi_1 I & \nabla c(x_k) \\ \nabla c(x_k)^T & -\xi_2 I \end{bmatrix}$$

where  $\xi_1$  **convexifies** the model and  $\xi_2$  **regularizes** the constraints

- ▶ Poor choices of  $\xi_1$  and  $\xi_2$  can have terrible consequences in the algorithm

1. *Journal of Management Studies*, 1990, 27, 1, 1-14.



- In computation of  $d_k = v_k + u_k$ , **convexify** the Hessian as in

$$\begin{bmatrix} H(x_k, \lambda_k) + \xi_1 I & \nabla c(x_k) \\ \nabla c(x_k)^T & 0 \end{bmatrix}$$

by monitoring iterates

- ▶ Hessian modification strategy: Increase  $\xi_1$  whenever

$$\begin{aligned} \|u_k\|^2 &> \psi \|v_k\|^2, \quad \psi > 0 \\ \frac{1}{2} u_k^T (H(x_k, \lambda_k) + \xi_1 I) u_k &< \theta \|u_k\|^2, \quad \theta > 0 \end{aligned}$$



## Algorithm 2: Inexact Newton (regularized)

(Curtis, Nocedal, Wächter (2009))

for  $k = 0, 1, 2, \dots$

- Approximately solve

$$\min \frac{1}{2} \|c(x_k) + \nabla c(x_k)^T v\|^2, \quad \text{s.t. } \|v\| \leq \omega \|(\nabla c(x_k))c(x_k)\|$$

to compute  $v_k$  satisfying **Cauchy decrease**

- Iteratively solve

$$\begin{bmatrix} H(x_k, \lambda_k) + \xi_1 I & \nabla c(x_k) \\ \nabla c(x_k)^T & 0 \end{bmatrix} \begin{bmatrix} d_k \\ \delta_k \end{bmatrix} = - \begin{bmatrix} \nabla f(x_k) + \nabla c(x_k) \lambda_k \\ -\nabla c(x_k)^T v_k \end{bmatrix}$$

**until termination test 1 or 2 is satisfied, increasing  $\xi_1$  as described**

- If only termination test 2 is satisfied, **increase  $\pi$**  so

$$\pi_k \geq \max \left\{ \pi_{k-1}, \frac{\nabla f(x_k)^T d_k + \max\{\frac{1}{2} u_k^T (H(x_k, \lambda_k) + \xi_1 I) u_k, \theta \|u_k\|^2\}}{(1 - \tau)(\|c(x_k)\| - \|c(x_k) + \nabla c(x_k)^T d_k\|)} \right\}$$

- Backtrack from  $\alpha_k \leftarrow 1$  to satisfy

$$\phi(x_k + \alpha_k d_k; \pi_k) \leq \phi(x_k; \pi_k) - \eta \alpha_k \Delta m(d_k; \pi_k)$$

- Update iterate  $(x_{k+1}, \lambda_{k+1}) \leftarrow (x_k, \lambda_k) + \alpha_k(d_k, \delta_k)$

## Assumption

The sequence  $\{(x_k, \lambda_k)\}$  is contained in a convex set  $\Omega$  over which  $f$ ,  $c$ , and their first derivatives are bounded and Lipschitz continuous

## Theorem

(Curtis, Nocedal, Wächter (2009)) If all limit points of  $\{\nabla c(x_k)^T\}$  have full row rank, then the sequence  $\{(x_k, \lambda_k)\}$  yields the limit

$$\lim_{k \rightarrow \infty} \left\| \begin{bmatrix} \nabla f(x_k) + \nabla c(x_k) \lambda_k \\ c(x_k) \end{bmatrix} \right\| = 0.$$

*Otherwise.*

$$\lim_{k \rightarrow \infty} \|(\nabla c(x_k))c(x_k)\| = 0$$

and if  $\{\pi_k\}$  is bounded, then

$$\lim_{k \rightarrow \infty} \|\nabla f(x_k) + \nabla c(x_k)\lambda_k\| = 0$$

# Handling inequalities

- ▶ **Interior point methods** are attractive for large applications
- ▶ Line-search interior point methods that enforce

$$c(x_k) + \nabla c(x_k)^T d_k = 0$$

may fail to converge globally (Wächter, Biegler (2000))

- ▶ Fortunately, the trust region subproblem we use to regularize the constraints also saves us from this type of failure!

## Algorithm 2 (Interior-point version)

- Apply Algorithm 2 to the logarithmic-barrier subproblem

$$\min f(x) - \mu \sum_{i=1}^q \ln s^i, \quad \text{s.t. } c_{\mathcal{E}}(x) = 0, \quad c_{\mathcal{I}}(x) - s = 0$$

for  $\mu \rightarrow 0$

- Define

$$\begin{bmatrix} H(x_k, \lambda_{\mathcal{E},k}, \lambda_{\mathcal{I},k}) & 0 & \nabla c_{\mathcal{E}}(x_k) & \nabla c_{\mathcal{I}}(x_k) \\ 0 & \mu I & 0 & -S_k \\ \nabla c_{\mathcal{E}}(x_k)^T & 0 & 0 & 0 \\ \nabla c_{\mathcal{I}}(x_k)^T & -S_k & 0 & 0 \end{bmatrix} \begin{bmatrix} d_k^x \\ d_k^s \\ \delta_{\mathcal{E},k} \\ \delta_{\mathcal{I},k} \end{bmatrix}$$

so that the iterate update has

$$\begin{bmatrix} x_{k+1} \\ s_{k+1} \end{bmatrix} \leftarrow \begin{bmatrix} x_k \\ s_k \end{bmatrix} + \alpha_k \begin{bmatrix} d_k^x \\ S_k d_k^s \end{bmatrix}$$

- Incorporate a fraction-to-the-boundary rule in the line search and a **slack reset** in the algorithm to maintain  $s \geq \max\{0, c_{\mathcal{I}}(x)\}$

# Convergence of Algorithm 2 (interior-point)

## Assumption

*The sequence  $\{(x_k, \lambda_{\varepsilon,k}, \lambda_{\mathcal{I},k})\}$  is contained in a convex set  $\Omega$  over which  $f$ ,  $c_{\varepsilon}$ ,  $c_{\mathcal{I}}$ , and their first derivatives are bounded and Lipschitz continuous*

## Theorem

*(Curtis, Schenk, Wächter (2009))*

- ▶ *For a given  $\mu$ , Algorithm 2 yields the same limits as in the equality constrained case*
- ▶ *If Algorithm 2 yields a sufficiently accurate solution to the barrier subproblem for each  $\{\mu_j\} \rightarrow 0$  and if the linear independence constraint qualification (LICQ) holds at a limit point  $\bar{x}$  of  $\{x_j\}$ , then there exist Lagrange multipliers  $\bar{\lambda}$  such that the first-order optimality conditions of the nonlinear program are satisfied*

# Outline

## PDE-Constrained Optimization

Introduction

Newton's method

Inexactness

**Results**

Summary and future work

## Nonsmooth Optimization

Sequential quadratic programming (SQP)

Gradient sampling (GS)

SQP-GS

Results

Summary and future work

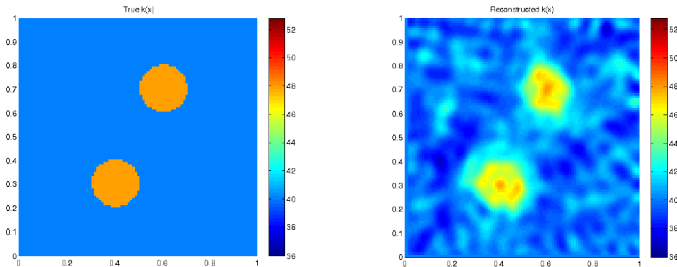
## Conclusion

# Implementation details

- ▶ Incorporated in IPOPT software package (Wächter)
  - ▶ `inexact_algorithm` yes
- ▶ Linear systems solved with PARDISO (Schenk)
  - ▶ SQMR (Freund (1994))
- ▶ Preconditioning in PARDISO
  - ▶ incomplete multilevel factorization with inverse-based pivoting
  - ▶ stabilized by symmetric-weighted matchings
- ▶ Optimality tolerance:  $1e-8$

- ▶ 745 problems written in AMPL
- ▶ 645 solved successfully
- ▶ 42 “real” failures
- ▶ Robustness between 87%-94%
- ▶ Original IPOPT: 93%





$N$	$n$	$p$	$q$	# iter	CPU sec (per iter)
32	14724	13824	1800	37	807.823 (21.833)
64	56860	53016	7688	25	3741.42 (149.66)
128	227940	212064	31752	20	54581.8 (2729.1)

# Boundary control

$$\begin{aligned}
 & \min \frac{1}{2} \int_{\Omega} (y(x) - y_t(x))^2 dx \\
 & \text{s.t. } -\nabla \cdot (e^{y(x)} \cdot \nabla y(x)) = 20 \quad \text{in } \Omega \\
 & \quad y(x) = u(x) \quad \text{on } \partial\Omega \\
 & \quad 2.5 \leq u(x) \leq 3.5 \quad \text{on } \partial\Omega
 \end{aligned}$$

where

$$y_t(x) = 3 + 10x_1(x_1 - 1)x_2(x_2 - 1)\sin(2\pi x_3)$$

$N$	$n$	$p$	$q$	# iter	CPU sec (per iter)
16	4096	2744	2704	13	2.8144 (0.2165)
32	32768	27000	11536	13	103.65 (7.9731)
64	262144	238328	47632	14	5332.3 (380.88)

Original IPOPT with  $N = 32$  requires 238 seconds per iteration

# Hyperthermia treatment planning

$$\begin{aligned}
 \min \quad & \frac{1}{2} \int_{\Omega} (y(x) - y_t(x))^2 dx \\
 \text{s.t.} \quad & -\Delta y(x) - 10(y(x) - 37) = u^* M(x) u \quad \text{in } \Omega \\
 & 37.0 \leq y(x) \leq 37.5 \quad \text{on } \partial\Omega \\
 & 42.0 \leq y(x) \leq 44.0 \quad \text{in } \Omega_0
 \end{aligned}$$

where

$$u_j = a_j e^{i\phi_j}, \quad M_{jk}(x) = \langle E_j(x), E_k(x) \rangle, \quad E_j = \sin(jx_1 x_2 x_3 \pi)$$

$N$	$n$	$p$	$q$	# iter	CPU sec (per iter)
16	4116	2744	2994	68	22.893 (0.3367)
32	32788	27000	13034	51	3055.9 (59.920)

Original IPOPT with  $N = 32$  requires 408 seconds per iteration

# Groundwater modeling

$$\begin{aligned} \min \quad & \frac{1}{2} \int_{\Omega} (y(x) - y_t(x))^2 dx + \frac{1}{2} \alpha \int_{\Omega} [\beta(u(x) - u_t(x))^2 + |\nabla(u(x) - u_t(x))|^2] dx \\ \text{s.t.} \quad & -\nabla \cdot (e^{u(x)} \cdot \nabla y_i(x)) = q_i(x) \quad \text{in } \Omega, \quad i = 1, \dots, 6 \\ & \nabla y_i(x) \cdot n = 0 \quad \text{on } \partial\Omega \\ & \int_{\Omega} y_i(x) dx = 0, \quad i = 1, \dots, 6 \\ & -1 \leq u(x) \leq 2 \quad \text{in } \Omega \end{aligned}$$

where

$$q_i = 100 \sin(2\pi x_1) \sin(2\pi x_2) \sin(2\pi x_3)$$

$N$	$n$	$p$	$q$	# iter	CPU sec (per iter)
16	28672	24576	8192	18	206.416 (11.4676)
32	229376	196608	65536	20	1963.64 (98.1820)
64	1835008	1572864	524288	21	134418. (6400.85)

Original IPOPT with  $N = 32$  requires approx. 20 **hours** for the first iteration

# Outline

## PDE-Constrained Optimization

Introduction

Newton's method

Inexactness

Results

Summary and future work

## Nonsmooth Optimization

Sequential quadratic programming (SQP)

Gradient sampling (GS)

SQP-GS

Results

Summary and future work

## Conclusion

# Summary

- ▶ We have a new framework for inexact Newton methods for optimization
- ▶ Convergence results are as good (and sometimes better) than exact methods
- ▶ Preliminary numerical results are encouraging

## Future work

- ▶ Tune the method for specific applications
- ▶ Incorporate useful techniques such as filters, second-order corrections, specialized preconditioners
- ▶ Use (approximate) elimination techniques so that larger (e.g., time-dependent) problems can be solved

# Outline

## PDE-Constrained Optimization

Introduction

Newton's method

Inexactness

Results

Summary and future work

## Nonsmooth Optimization

Sequential quadratic programming (SQP)

Gradient sampling (GS)

SQP-GS

Results

Summary and future work

## Conclusion



# Constrained optimization of smooth functions

- Consider constrained optimization problems of the form

$$\begin{aligned} \min_x \quad & f(x) \\ \text{s.t.} \quad & c(x) \leq 0 \end{aligned}$$

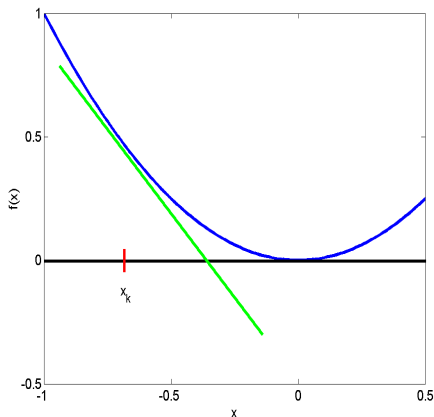
where  $f$  and  $c$  are *smooth* (equality constraints OK, too)

- At  $x_k$ , solve the SLP/SQP subproblem

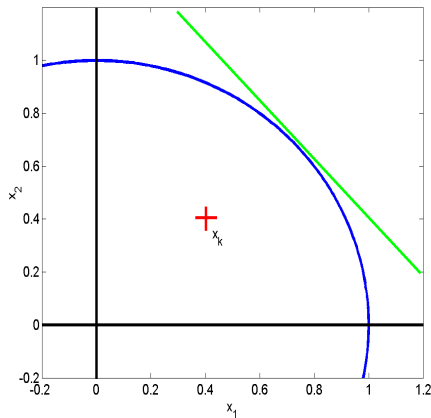
$$\begin{aligned} \min_d \quad & f_k + \nabla f_k^T d + \frac{1}{2} d^T H_k d \\ \text{s.t.} \quad & c_k + \nabla c_k^T d \leq 0, \quad \|d\| \leq \Delta_k \end{aligned}$$

to compute the search direction  $d_k$

# SQP illustration: Objective model



# SQP illustration: Constraint model



# Practicalities

- ▶ Since the linearized constraints may be inconsistent, we solve

$$\begin{aligned} \min_d \quad & \rho(f_k + \nabla f_k^T d) + \sum s^i + \frac{1}{2} d^T H_k d \\ \text{s.t.} \quad & c_k + \nabla c_k^T d \leq s, \quad s \geq 0, \end{aligned}$$

where  $\rho > 0$  is a *penalty parameter*

- ▶ We perform a line search on the penalty function

$$\phi(x; \rho) \triangleq \rho f(x) + \sum \max\{0, c^i(x)\}$$

to promote global convergence

# Line search

- ▶ A model of the penalty function is given by

$$q_k(d; \rho) \triangleq \rho(f_k + \nabla f_k^T d) + \sum \max\{0, c_k^i + \nabla c_k^{i^T} d\} + \frac{1}{2} d^T H_k d$$

- ▶ Solving the SQP subproblem is equivalent to minimizing  $q_k(d; \rho)$
- ▶ The reduction in  $q_k(d; \rho)$  yielded by  $d_k$  is

$$\Delta q_k(d_k; \rho) \triangleq q_k(0; \rho) - q_k(d_k; \rho)$$

- ▶ We impose the sufficient decrease condition

$$\phi(x_k + \alpha_k d_k; \rho) \leq \phi(x_k; \rho) - \eta \alpha_k \Delta q_k(d_k; \rho)$$

## Penalty-SQP method

for  $k = 0, 1, 2, \dots$

- Solve the SQP subproblem

$$\begin{aligned} \min_d \quad & \rho(f_k + \nabla f_k^T d) + \sum s^i + \frac{1}{2} d^T H_k d \\ \text{s.t.} \quad & c_k + \nabla c_k^T d \leq s, \quad s \geq 0 \end{aligned}$$

or, equivalently, solve

$$\min_d q_k(d; \rho) \triangleq \rho(f_k + \nabla f_k^T d) + \sum \max\{0, c_k^i + \nabla c_k^{iT} d\} + \frac{1}{2} d^T H_k d$$

to compute  $d_k$

- Backtrack from  $\alpha_k = 1$  to satisfy

$$\phi(x_k + \alpha_k d_k; \rho) \leq \phi(x_k; \rho) - \eta \alpha_k \Delta q_k(d_k; \rho)$$

- Update  $x_{k+1} \leftarrow x_k + \alpha_k d_k$

# Outline

## PDE-Constrained Optimization

Introduction

Newton's method

Inexactness

Results

Summary and future work

## Nonsmooth Optimization

Sequential quadratic programming (SQP)

**Gradient sampling (GS)**

SQP-GS

Results

Summary and future work

## Conclusion

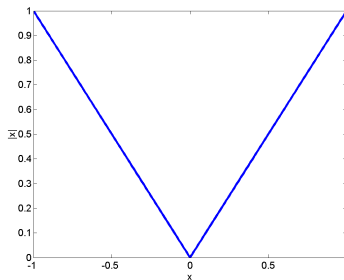
# Unconstrained optimization of nonsmooth functions

- Consider the unconstrained optimization problem

$$\min_x f(x)$$

where  $f$  may be nonsmooth (but is at least locally Lipschitz)

- The prototypical example is the absolute value function:





# Clarke subdifferential

- ▶ Suppose  $f$  is differentiable over an open dense set  $\mathcal{D}$
- ▶ Let

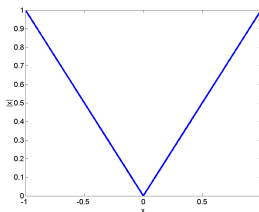
$$\mathbb{B}(x', \epsilon) \triangleq \{x \mid \|x - x'\| \leq \epsilon\}$$

- ▶ The Clarke subdifferential is

$$\bar{\partial}f(x') = \bigcap_{\epsilon > 0} \text{cl conv } \nabla f(\mathbb{B}(x', \epsilon) \cap \mathcal{D})$$

- ▶ A point  $x'$  is called Clarke stationary if  $0 \in \bar{\partial}f(x')$

- ▶ A point  $x'$  is called  $\epsilon$ -stationary if  $0 \in \bar{\partial}f(x', \epsilon)$



- ... find  $\epsilon$ -stationary point, reduce  $\epsilon$ , find  $\epsilon$ -stationary point,...

# Gradient sampling: Robust steepest descent

- ▶ (Burke, Lewis, Overton, 2005)
- ▶ We restrict iterates to the open dense set  $\mathcal{D}$
- ▶ Ideally, at  $x_k$ , for a given  $\epsilon$  we would solve

$$\min_d f_k + \max_{x \in \mathcal{B}_k} \{\nabla f(x)^T d\} + \frac{1}{2} d^T H_k d$$

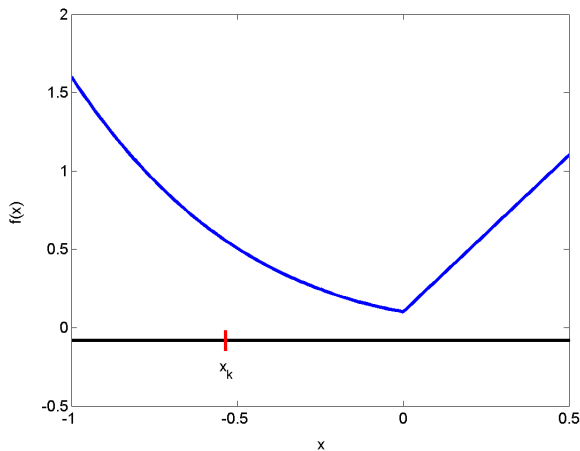
where  $\mathcal{B}_k = \mathbb{B}(x_k, \epsilon) \cap \mathcal{D}$

- ▶ However, we can only approximate this step by solving

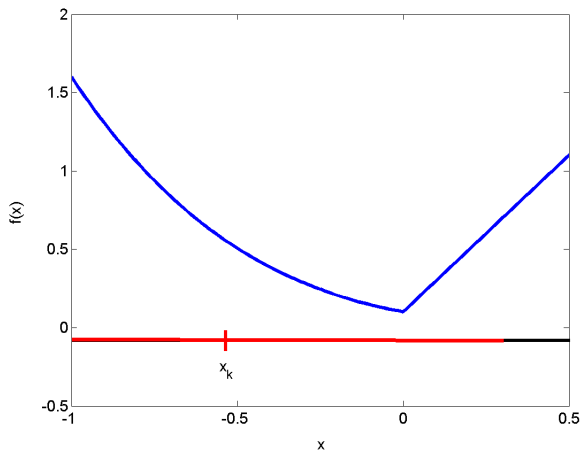
$$\min_d f_k + \max_{x \in \mathcal{B}_k} \{\nabla f(x)^T d\} + \frac{1}{2} d^T H_k d$$

where  $\mathcal{B}_k = \{x_k, x_{k1}, \dots, x_{kp}\} \subset \mathbb{B}(x_k, \epsilon) \cap \mathcal{D}$

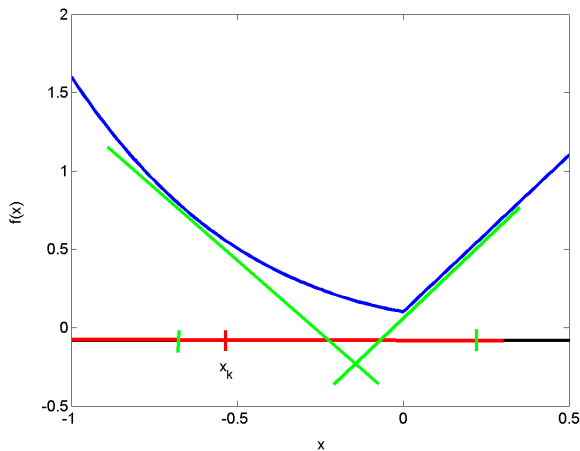
# GS illustration: Objective model



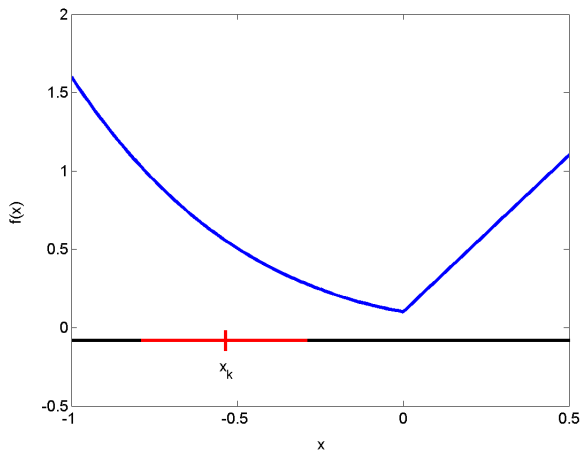
# GS illustration: Objective model



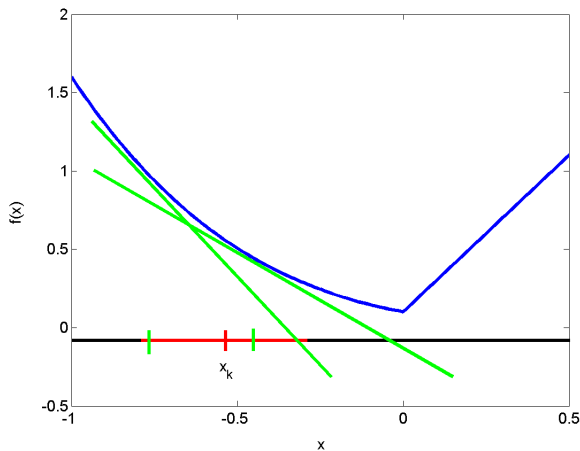
# GS illustration: Objective model



# GS illustration: Objective model



# GS illustration: Objective model





## GS method

for  $k = 0, 1, 2, \dots$

- ▶ Sample points  $\{x_{k1}, \dots, x_{kp}\}$  in  $\mathbb{B}(x_k, \epsilon) \cap \mathcal{D}$
- ▶ Solve the GS subproblem

$$\min_d f_k + \max_{x \in \mathcal{B}_k} \{\nabla f(x)^T d\} + \frac{1}{2} d^T H_k d$$

to compute  $d_k$

- ▶ Backtrack from  $\alpha_k = 1$  to satisfy

$$f(x_k + \alpha_k d_k) \leq f(x_k) - \eta \alpha_k \|d_k\|^2$$

- ▶ Update  $x_{k+1} \approx x_k + \alpha_k d_k$  (to ensure  $x_{k+1} \in \mathcal{D}$ )
- ▶ If  $\|d_k\| \leq \epsilon$ , then reduce  $\epsilon$

# Outline

## PDE-Constrained Optimization

Introduction

Newton's method

Inexactness

Results

Summary and future work

## Nonsmooth Optimization

Sequential quadratic programming (SQP)

Gradient sampling (GS)

**SQP-GS**

Results

Summary and future work

## Conclusion

# Constrained optimization of nonsmooth functions

- Consider constrained optimization problems of the form

$$\begin{aligned} \min_x \quad & f(x) \\ \text{s.t.} \quad & c(x) \leq 0 \end{aligned}$$

where  $f$  and  $c$  may be *nonsmooth* (equality constraints OK, too)

- We may consider solving

$$\min_x \phi(x; \rho) \triangleq \rho f(x) + \sum \max\{0, c^i(x)\}$$

or even

$$\min_x \varphi(x; \rho) \triangleq \rho f(x) + \max_i \max\{0, c^i(x)\}$$

but this makes me... :-)

# SQP and GS

- The SQP subproblem is

$$\min_d \rho z + \sum s^i + \frac{1}{2} d^T H_k d$$

$$\text{s.t. } f_k + \nabla f_k^T d \leq z$$

$$c_k + \nabla c_k^T d \leq s, \quad s \geq 0$$

- The GS subproblem is

$$\min_d z + \frac{1}{2} d^T H_k d$$

$$\text{s.t. } f_k + \nabla f(x)^T d \leq z, \quad \forall x \in \mathcal{B}_k$$

# SQP-GS

- The SQP-GS subproblem is

$$\begin{aligned} \min_{d,z,s} \quad & \rho z + \sum s^i + \frac{1}{2} d^T H_k d \\ \text{s.t.} \quad & f_k + \nabla f(x)^T d \leq z, \quad \forall x \in \mathcal{B}_k^0 \\ & c_k^i + \nabla c^i(x)^T d \leq s^i, \quad s^i \geq 0, \quad \forall x \in \mathcal{B}_k^i, \quad i = 1, \dots, m \end{aligned}$$

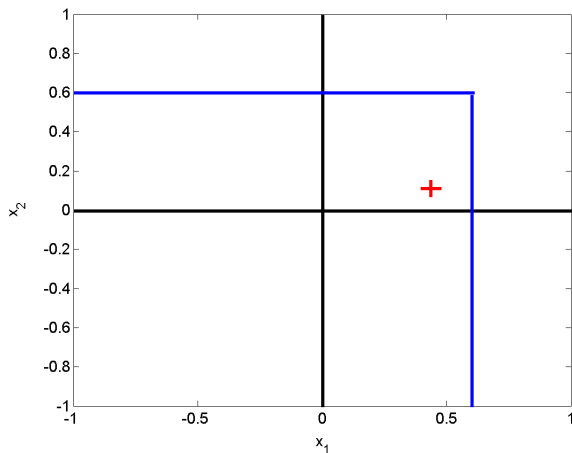
where  $\mathcal{B}_k^i = \{x_k, x_{k1}^i, \dots, x_{kp}^i\} \subset \mathbb{B}(x_k, \epsilon)$  for  $i = 0, \dots, m$

- This is equivalent to

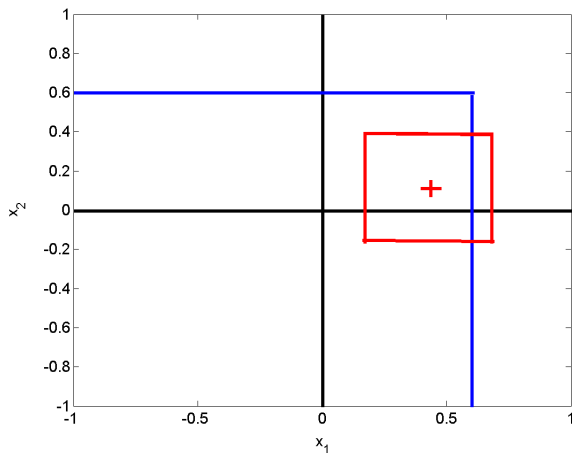
$$\min_d \quad \rho \max_{x \in \mathcal{B}_k^0} (f_k + \nabla f(x)^T d) + \sum_{x \in \mathcal{B}_k^i} \max \{0, c_k^i + \nabla c^i(x)^T d, 0\} + \frac{1}{2} d^T H_k d$$

i.e.,  $\min_d q_k(d; \rho)$ , where now  $q_k(d; \rho)$  is a *robust* model of  $\phi(x; \rho)$

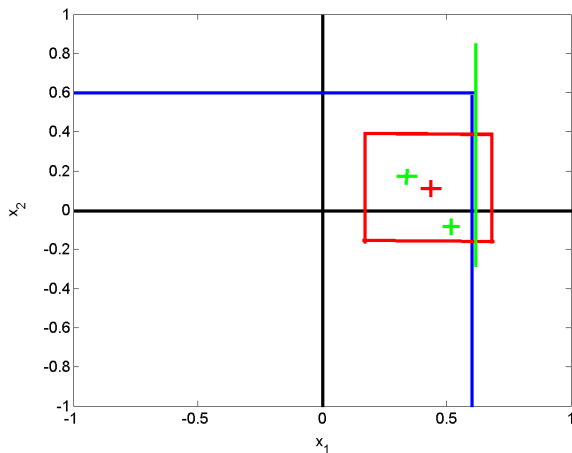
# SQP-GS illustration: Constraint model



# SQP-GS illustration: Constraint model

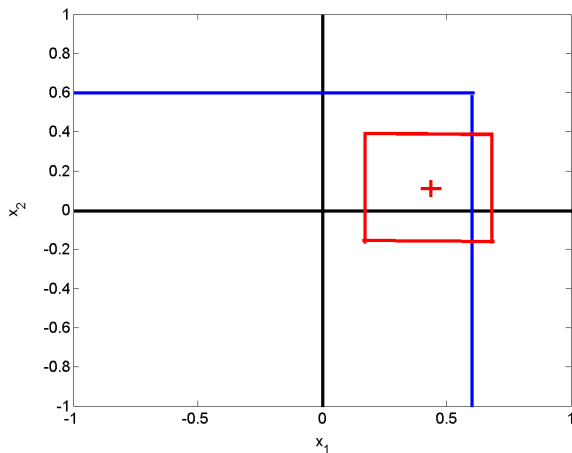


# SQP-GS illustration: Constraint model

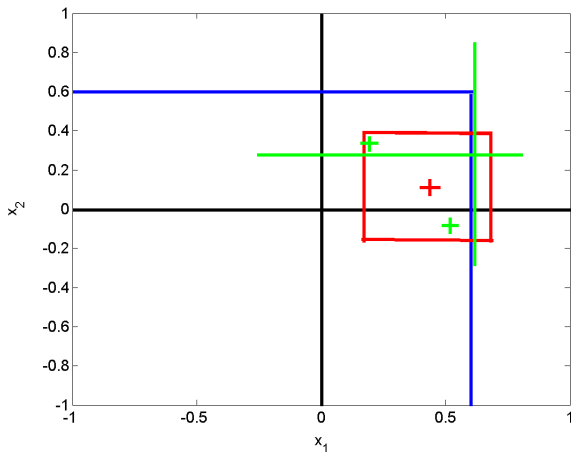




# SQP-GS illustration: Constraint model



# SQP-GS illustration: Constraint model



## SQP-GS method

for  $k = 0, 1, 2, \dots$

- ▶ Sample points  $\{x_{k1}^i, \dots, x_{kp}^i\}$  in  $\mathbb{B}(x_k, \epsilon) \in \mathcal{D}^i$  for  $i = 0, \dots, m$
- ▶ Solve the SQP-GS subproblem

$$\begin{aligned} \min_{d, z, s} \quad & \rho z + \sum s^i + \frac{1}{2} d^T H_k d \\ \text{s.t.} \quad & f_k + \nabla f(x)^T d \leq z, \quad \forall x \in \mathcal{B}_k^0 \\ & c_k^i + \nabla c^i(x)^T d \leq s^i, \quad s^i \geq 0, \quad \forall x \in \mathcal{B}_k^i, \quad i = 1, \dots, m \end{aligned}$$

to compute  $d_k$

- ▶ Backtrack from  $\alpha_k = 1$  to satisfy

$$\phi(x_k + \alpha_k d_k; \rho) \leq \phi(x_k; \rho) - \eta \alpha_k \Delta q_k(d_k; \rho)$$

- ▶ Update  $x_{k+1} \approx x_k + \alpha_k d_k$  (to ensure  $x_{k+1} \in \cap_i \mathcal{D}^i$ )
- ▶ If  $\Delta q_k(d_k; \rho) \leq \epsilon$ , then reduce  $\epsilon$

# Global convergence

- ▶ Assumption 1: The functions  $f$  and  $c^i$ ,  $i = 1, \dots, m$ , are locally Lipschitz and continuously differentiable on open dense subsets of  $\mathbb{R}^n$
- ▶ Assumption 2: The sequence of iterates and sample points are contained in a convex set over which the functions  $f$  and  $c^i$ ,  $i = 1, \dots, m$ , and their first derivatives are bounded
- ▶ Assumption 3: For universal constants  $\bar{\xi} \geq \underline{\xi} > 0$ , the Hessian matrices satisfy  $\underline{\xi}\|d\|^2 \leq d^T H_k d \leq \bar{\xi}\|d\|^2$  for all  $d \in \mathbb{R}^n$

## Global convergence

- ▶ Lemma 1:  $\Delta q_k(d_k; \rho) = 0$  if and only if  $x_k$  is  $\epsilon$ -stationary
- ▶ Lemma 2: The one-sided directional derivative of the penalty function satisfies

$$\phi'(d_k; \rho) \leq d_k^T H_k d_k < 0$$

and so  $d_k$  is a descent direction for  $\phi(x; \rho)$  at  $x_k$

- ▶ **Lemma 3:** Suppose the sample size is  $p \geq n + 1$ . If the current iterate  $x_k$  is sufficiently close to a stationary point  $x'$  of the penalty function  $\phi(x; \rho)$ , then there exists a nonempty open set of sample sets such that the solution to the SQP-GS subproblem  $d_k$  yields an arbitrarily small  $\Delta q_k(d_k; \rho)$ 
  - ▶ Carathéodory's Theorem
- ▶ Theorem: With probability one, every cluster point of  $\{x_k\}$  is feasible and stationary for  $\phi(x; \rho)$

# Outline

## PDE-Constrained Optimization

Introduction

Newton's method

Inexactness

Results

Summary and future work

## Nonsmooth Optimization

Sequential quadratic programming (SQP)

Gradient sampling (GS)

SQP-GS

Results

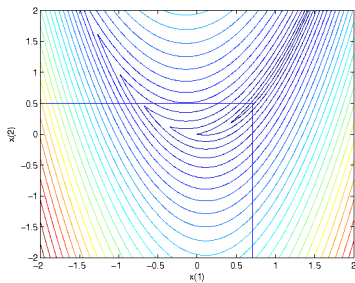
Summary and future work

## Conclusion

# Implementation

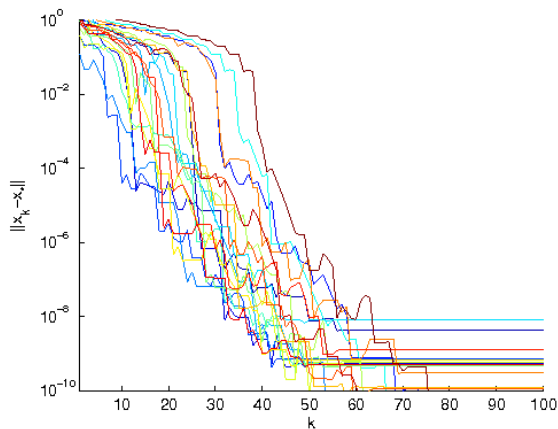
- ▶ Prototype implementation in MATLAB (available soon?)
- ▶ QP subproblems solved with MOSEK
- ▶ BFGS approximations of Hessian of penalty function
  - ▶ (Lewis and Overton, 2009)
- ▶  $\rho$  decreased conservatively

$$\begin{aligned} \min_x \quad & 8|x_1^2 - x_2| + (1 - x_1)^2 \\ \text{s.t.} \quad & \max\{\sqrt{2}x_1, 2x_2\} \leq 1 \end{aligned}$$





# Example 1: Nonsmooth Rosenbrock



## Example 2: Entropy minimization

Find a  $N \times N$  matrix  $X$  that solves

$$\begin{aligned} \min_X \quad & \ln \left( \prod_{j=1}^K \lambda_j(A \circ X^T X) \right) \\ \text{s.t.} \quad & \|X_j\| = 1, \quad j = 1, \dots, N \end{aligned}$$

where  $\lambda_j(M)$  denotes the  $j$ th largest eigenvalue of  $M$ ,  $A$  is a real symmetric  $N \times N$  matrix,  $\circ$  denotes the Hadamard matrix product, and  $X_j$  denotes the  $j$ th column of  $X$

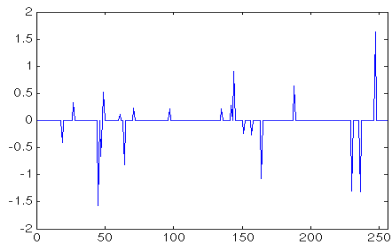
$N$	$n$	$K$	$f$ (SQP-GS)	$f$ (GS)
2	4	1	1.00000e+00	1.00000e+00
4	16	2	7.46296e-01	7.46286e-01
6	36	3	6.33589e-01	6.33477e-01
8	64	4	5.60165e-01	5.58820e-01
10	100	5	2.20724e-01	2.17193e-01
12	144	6	1.24820e-01	1.22226e-01
14	196	7	8.21835e-02	8.01010e-02
16	256	8	5.73762e-02	5.57912e-02

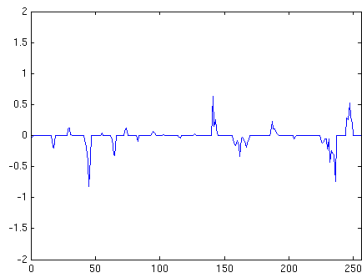
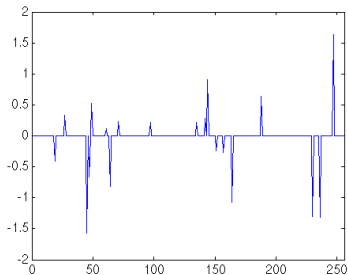
## Example 3(a): Compressed sensing ( $\ell_1$ norm)

Recover a sparse signal by solving

$$\begin{aligned} \min_x \quad & \|x\|_1 \\ \text{s.t.} \quad & Ax = b \end{aligned}$$

where  $A$  is a  $64 \times 256$  submatrix of a discrete cosine transform (DCT) matrix



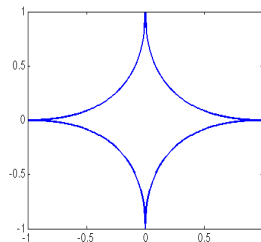
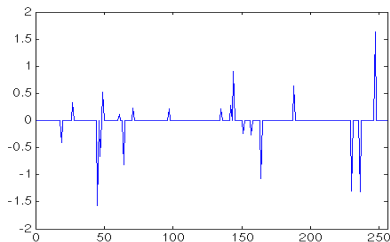


# Example 3(b): Compressed sensing ( $\ell_{0.5}$ norm)

Recover a sparse signal by solving

$$\begin{aligned} \min_x \quad & \|x\|_{0.5} \\ \text{s.t.} \quad & Ax = b \end{aligned}$$

where  $A$  is a  $64 \times 256$  submatrix of a discrete cosine transform (DCT) matrix



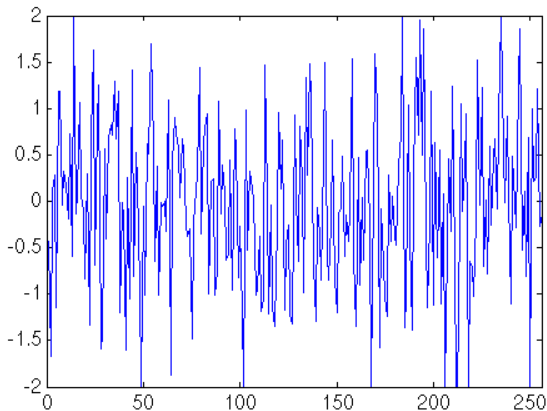


Figure:  $k = 1$

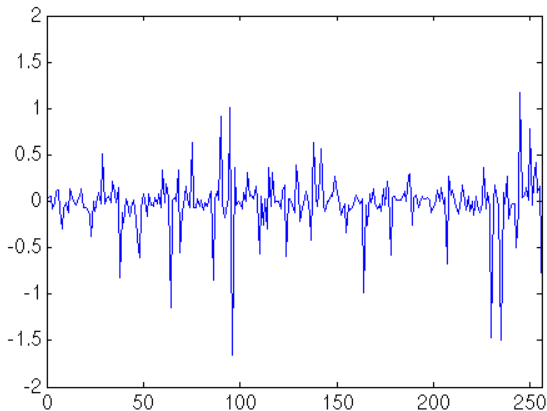


Figure:  $k = 10$



1.  $\frac{1}{2} \times \frac{1}{2} = \frac{1}{4}$

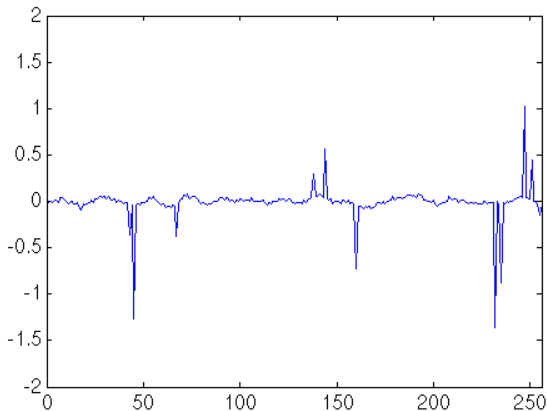


Figure:  $k = 25$

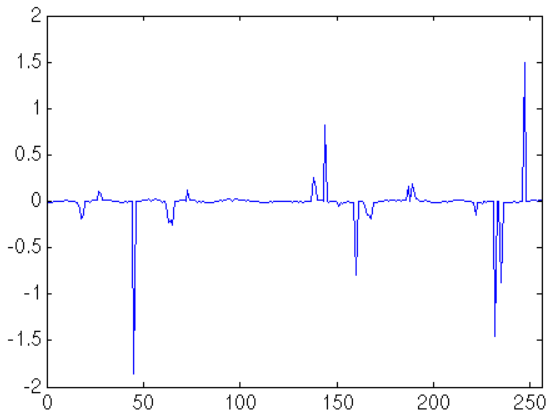


Figure:  $k = 50$

1.  $\frac{1}{2} \times \frac{1}{2} = \frac{1}{4}$

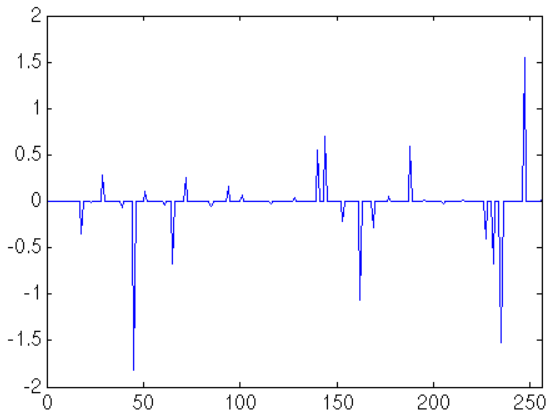
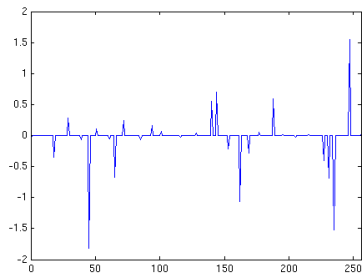
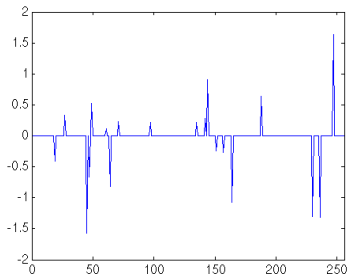


Figure:  $k = 200$



# Outline

## PDE-Constrained Optimization

Introduction

Newton's method

Inexactness

Results

Summary and future work

## Nonsmooth Optimization

Sequential quadratic programming (SQP)

Gradient sampling (GS)

SQP-GS

Results

Summary and future work

## Conclusion

# Summary

- ▶ We have presented a globally convergent algorithm for the solution of constrained, nonsmooth, and nonconvex optimization problems
- ▶ The algorithm follows a penalty-SQP framework and uses Gradient Sampling to make the search direction calculation robust
- ▶ Preliminary results are encouraging

## Future work

- ▶ Tune updates for  $\epsilon$  and  $\rho$
- ▶ Allow for special handling of smooth/convex/linear functions
- ▶ Investigate SLP vs. SQP
- ▶ Extensions for particular applications; e.g., specialized sampling

Thanks!!