

$$\phi''(0) = \lim_{t \rightarrow 0} \frac{\phi(t) - \phi(0)}{t^2} < 0$$

Pick t small enough $\frac{\phi(t) - \phi(0)}{t^2} \leq \frac{\phi''(0)}{2}$.

$$\Rightarrow \phi(t) \leq \phi(0) + \underbrace{t^2 \frac{\phi''(0)}{2}}_{< 0}$$

$$\Rightarrow f(x^* + ts) < f(x^*)$$

\Rightarrow Not local min. □

Starting 9/7

Last time, we showed that necessarily every local minimum of some $f(x)$ has

$$\nabla f(x^*) = 0, \quad \nabla^2 f(x^*) \geq 0.$$

"zero gradient" "positive semidefinite Hessian"

If $\nabla f(x^*) \neq 0$, then $x^* \leftarrow x^* - t \nabla f(x^*)$

If $\nabla^2 f(x^*) \not\geq 0$, then $x^* \leftarrow x^* - t s$

\Downarrow
 \exists negative eigenvalue

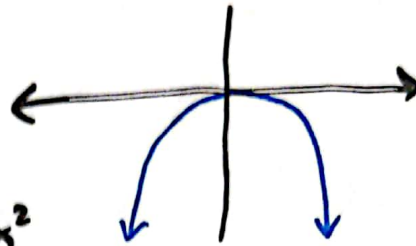
\nwarrow most negative eigenvector in $\nabla^2 f(x^*)$

Is it sufficient to have zero grad and p.s.d. Hessian?

No! $f(x) = -x^4$

$$f'(0) = 0 \text{ since}$$

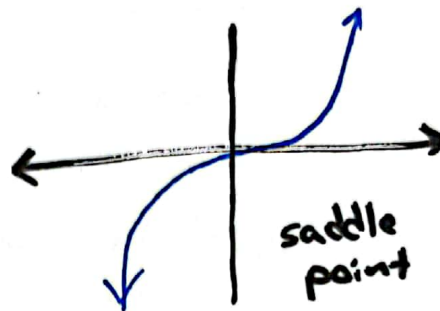
$$f'(x) = -4x^3$$



local max

$$f(x) = x^3$$

$$f'(0) = 0$$



saddle point

5. Theorem (Second-Order Sufficient Condition)

Suppose f is twice diff.

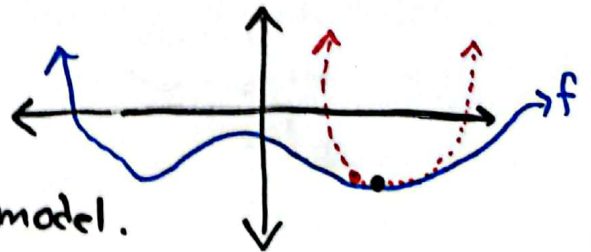
If $x^* \in \mathbb{R}^d$ satisfies $\nabla f(x^*) = 0$ and $\nabla^2 f(x^*) \succ 0$,



then x^* is a strict local minimum. $\forall s \neq 0, s^T \nabla^2 f(x^*) s > 0$.

Proof Idea. Look in 1D.

Strict local min
of second-order model.



Proof. Observe that $\lambda_{\min}(\nabla^2 f(\bar{x})) > 0$ (strictly).

For some $\varepsilon > 0$, $\lambda > 0$, every $x \in B(\bar{x}, \varepsilon)$

has $\lambda_{\min}(\nabla^2 f(x)) > \lambda > 0$,

(since $\nabla^2 f(x)$ is continuous and $\lambda_{\min}(\cdot)$ is a continuous.)

The Fundamental Theorem of Calculus

$$\theta(1) = \theta(0) + \int_0^1 \theta'(t) dt$$

and for second derivatives, we have

$$\theta(1) = \theta(0) + \theta'(0) + \int_0^1 \int_0^t \theta''(\alpha) d\alpha dt$$

Part of HW 1

For $s \in \mathbb{R}^d$, $\theta(t) = f(\bar{x} + ts)$, $\theta'(0) = \nabla f(\bar{x})^T s$
 $\|s\| \leq \varepsilon$ $= 0$

$$\begin{aligned} \theta''(\alpha) &= s^T \nabla^2 f(\bar{x} + \alpha s) s \\ &\geq \lambda \|s\|_2^2 > 0 \end{aligned}$$

For any s ,

$$\begin{aligned} f(\bar{x} + s) &\geq f(\bar{x}) + 0 + \lambda \|s\|_2^2 \underbrace{\int_0^1 \int_0^t 1 d\alpha dt}_{= 1/2} \\ &> f(\bar{x}) \end{aligned}$$

if $s \neq 0$.

\Rightarrow Every nearby point is strictly worse than \bar{x} . \square

Two Example Optimization Problems

Example 1 (Least Squares/ Data Fitting)

In 1801, Gauss predicts
the location Ceres (asteroid/
dwarf planet).

Given several past locations x

In particular, 22 observations

$$(x_1, y_1), \dots, (x_N, y_N), N=22.$$

Model this as an ellipse (conic section)

$$(*) \quad \alpha x^2 + \beta y^2 + \gamma xy = 1 \quad \text{for some } \alpha, \beta, \gamma$$



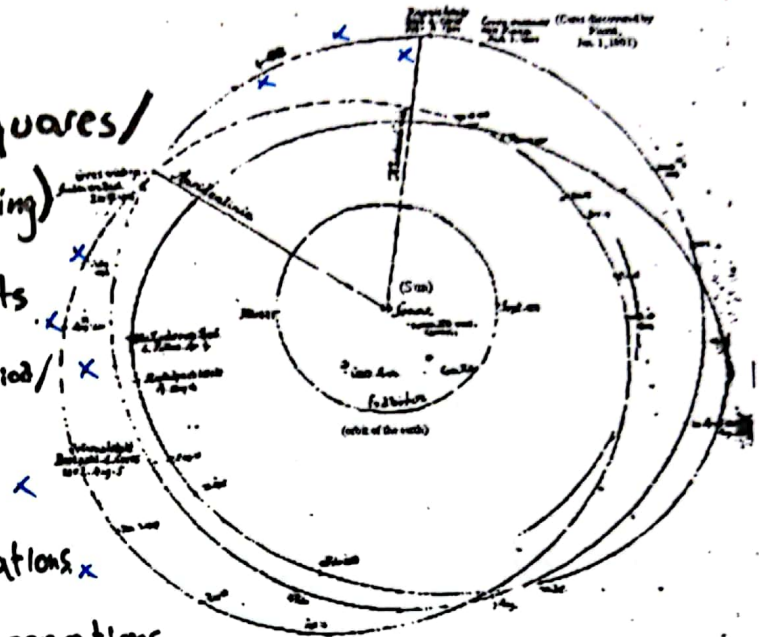
Gauss

What is the best fitting ellipse?

Minimize how much we violate $(*)$ squared:

$$\min_{(\alpha, \beta, \gamma) \in \mathbb{R}^3} \sum_{i=1}^N (\alpha x_i^2 + \beta y_i^2 + \gamma x_i y_i - 1)^2$$

Gauss solved this and succeeded



More generally, this is least squares optimization

$$\min_{x \in \mathbb{R}^d} \|Ax - b\|_2^2 = \sum_{i=1}^n (a_i^T x - b_i)^2$$

Selecting $x = \begin{bmatrix} \alpha \\ \beta \\ \gamma \end{bmatrix}$ $b_i = 1$, $a_i = \begin{bmatrix} x_i^2 \\ y_i^2 \\ x_i y_i \end{bmatrix}$

If we want more "accurate" model, we search for

$$\alpha x_i^4 + \beta y_i^4 + \gamma x_i^3 y_i + \dots + \lambda x_i + \omega y_i = 1$$

$$\Rightarrow x = \begin{bmatrix} \alpha \\ \beta \\ \gamma \\ \vdots \\ \omega \end{bmatrix}, b_i = 1, a_i = \begin{bmatrix} x_i^4 \\ y_i^4 \\ \vdots \\ x_i \\ y_i \end{bmatrix}.$$

HW2 will show $\|Ax - b\|_2^2$ is smooth, convex, and then solve it very effectively.

Example 2 (Logistic Regression)

Suppose you want "quick and dirty" test for if someone might have corona.

Given a patient, we know

c_1 = Patient's age

c_2 = blood pressure

c_3 = heart rate

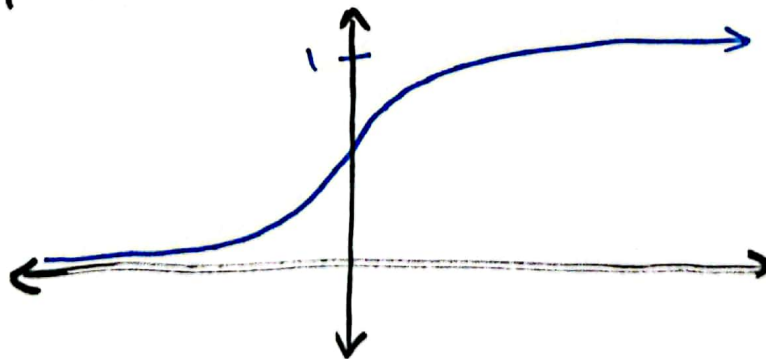
c_4 = temp

c_5 = Family history of cancer

⋮

Lets build a linear model $c^T x$ and use

$$\frac{e^{c^T x}}{1 + e^{c^T x}} \approx P(\text{patient has corona})$$



Given past patients c_i , and outcome

$$r_i = \begin{cases} 0 & \text{if did not have corona} \\ 1 & \text{otherwise.} \end{cases}$$

$$\min_{x \in \mathbb{R}^d} \left(\sum_{\substack{i \\ s.t. r_i = 0}} \frac{e^{c_i^T x}}{1 + e^{c_i^T x}} - \sum_{\substack{i \\ s.t. r_i = 1}} \frac{e^{c_i^T x}}{1 + e^{c_i^T x}} \right) \approx \sum_i \left| \frac{e^{c_i^T x}}{1 + e^{c_i^T x}} - r_i \right|$$

↑ smooth but not convex.

Neural Network model could replace

$$c^T x \leftarrow x \otimes \phi(\dots \phi(A_2 \phi(A_1 \cdot)) \cdot)$$

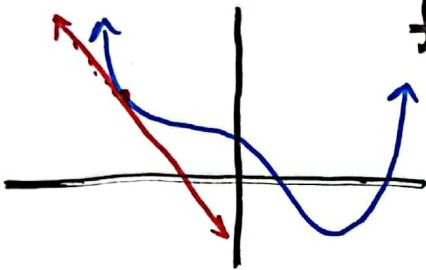
↑ elementwise $\max\{x, 0\}$

↑ nonsmooth
and nonconvex

First / Second - Order Iterative Improvement

Given $x_k \in \mathbb{R}^d$, how do we find a better point x_{k+1} given $(\nabla f(x_k), \nabla^2 f(x_k))$?

"First-Order Approach"



$$f(x_k + s) \approx f(x_k) + \nabla f(x_k)^T s$$

Attempt 1

$$\begin{aligned} \min_{s \in \mathbb{R}^d} & f(x_k) + \nabla f(x_k)^T s \\ = & \begin{cases} -\infty & \text{if } \nabla f(x_k) \neq 0 \\ f(x_k) & \text{otherwise.} \end{cases} \end{aligned}$$

Attempt 2

$$\begin{aligned} \min_{\|s\|_2 \leq \epsilon} & f(x_k) + \nabla f(x_k)^T s \end{aligned}$$

$$x_{k+1} \leftarrow x_k + \underset{\| \cdot \|}{\operatorname{argmin}} \left\{ f(x_k) + \nabla f(x_k)^T s \mid \|s\|_2 \leq \epsilon \right\}$$
$$= x_k - \frac{\nabla f(x_k)}{\|\nabla f(x_k)\|} \epsilon$$

Claim: $-\epsilon \frac{\nabla f(x_k)}{\|\nabla f(x_k)\|} = \operatorname{argmin}_s \{ \nabla f(x_k)^T s \mid \|s\|_2 \leq \epsilon \}$

Proof. Suppose for contradiction that s' does better

$$\begin{aligned} \nabla f(x_k)^T s' &< \nabla f(x_k)^T \left(\frac{-\epsilon \nabla f(x_k)}{\|\nabla f(x_k)\|_2} \right) \\ &= -\epsilon \|\nabla f(x_k)\|_2 \end{aligned}$$

By Cauchy-Schwartz,

$$-\|\nabla f(x_k)\|_2 \|s'\|_2 < -\epsilon \|\nabla f(x_k)\|_2$$

$$\Rightarrow \|s'\|_2 > \epsilon.$$

$$\Rightarrow s' \text{ is not feasible.} \quad \square$$

$$\Rightarrow x_{k+1} = x_k - \alpha_k \nabla f(x_k) \quad \text{for some } \alpha_k.$$

"Gradient Descent".