

# Nonlinear Optimization in Machine Learning

## Lecture 1 Introduction & Foundations

Why nonlinear optimization?

motivated by Machine Learning Applications

$$\mathcal{D} = \{(a_j, y_j), j=1, 2, \dots, m\}$$

"learn"  $\phi = \phi(a; x)$

$$\begin{aligned} \text{"loss"} \quad L_{\mathcal{D}}(x) &= \sum_{j=1}^m \mathcal{L}(a_j, y_j; x) \\ &= \sum_{j=1}^m \ell(\phi(a_j; x), y_j) \end{aligned}$$

$$x^* = \min_{x \in U} L_{\mathcal{D}}(x)$$

### Example 1 Least Squares

$$\min_x \frac{1}{2m} \sum_{j=1}^m (a_j^T x - y_j)^2 = \frac{1}{2m} \|Ax - y\|_2^2$$

"regularization"

$$\min_x \frac{1}{2m} \|Ax - y\|_2^2 + \lambda \|x\|_2^2 \quad (\lambda > 0)$$

(Tikhonov regularization)

$$\min_x \frac{1}{2m} \|Ax - y\|_2^2 + \lambda \|x\|_1$$

(LASSO: Least Absolute Shrinkage and Selection Operator)

Example 2 Matrix Completion

$A_j$  is  $n \times p$  and  $X$  is  $n \times p$

$$\min_x \frac{1}{2m} \sum_{j=1}^m (\langle A_j, X \rangle - y_j)^2$$

where  $\langle A, B \rangle = \text{tr}(A^T B)$

$$\min_x \frac{1}{2m} \sum_{j=1}^m (\langle A_j, X \rangle - y_j)^2 + \lambda \|X\|_*$$

$\|X\|_* = \text{sum of singular values of } X = \text{tr} \sqrt{X^T X}$   
= nuclear norm

$$L \in \mathbb{R}^{n \times r} \text{ and } R \in \mathbb{R}^{p \times r} \quad r \ll \min(n, p)$$

$$\min_{L, R} \frac{1}{2m} \sum_{j=1}^m (\langle A_j, LR^T \rangle - y_j)^2$$

Example 3 Nonnegative matrix factorization

$$\min_{L, R} \|LR^T - Y\|_F^2, \quad L \geq 0, R \geq 0$$

$$Y \in \mathbb{R}^{n \times p} \quad L \in \mathbb{R}^{n \times r} \quad R \in \mathbb{R}^{p \times r}$$

Example 4 Sparse inverse covariance estimation

$$\text{Sample covariance matrix} \quad S = \frac{1}{m-1} \sum_{j=1}^m a_j a_j^T$$

$$S^{-1} = X$$

$$\min_{X \in \text{Symmetric } \mathbb{R}^{n \times n}} \langle S, X \rangle - \log \det |X| + \lambda \|X\|_1$$

"Graphical Lasso"

$$X \geq 0$$

$$\|X\|_1 = \sum_{i, l=1}^n |X_{i, l}|$$

### Example 5 Sparse PCA

PCA = Principle Component analysis

$$\max_{V \in \mathbb{R}^n} V^T S V \quad \text{s.t.} \quad \|V\|_2 = 1, \|V\|_0 \leq k$$

"Sparse" via  $R \quad M = V V^T$

$$\max_{M \in \text{Symmetric } \mathbb{R}^{n \times n}} \langle S, M \rangle \quad \text{s.t.} \quad M \succeq 0, \langle I, M \rangle = 1, \|M\|_1 \leq R$$

Example 6 SVM (Support Vector Machine)

$$a_j \in \mathbb{R}^n \quad y_j \in \{-1, 1\}$$

$$\text{seek } x \in \mathbb{R}^n, \beta \in \mathbb{R} \quad \text{s.t.} \quad \begin{aligned} a_j^T x - \beta &\geq 1 \text{ if } y_j = 1 \\ a_j^T x - \beta &\leq -1 \text{ if } y_j = -1 \end{aligned}$$

$$H(x, \beta) = \frac{1}{m} \sum_{j=1}^m \max(1 - y_j(a_j^T x - \beta), 0)$$



# Example 7 Neural Network

"activation function"  $\sigma$

$$a_j^l = \sigma(W^l a_j^{l-1} + g^l), \quad l=1, 2, \dots, D$$

"weight"  $w = (W^1, g^1, W^2, g^2, \dots, W^D, g^D)$

$$L(w, X) = \frac{1}{m} \sum_{j=1}^m \left[ \sum_{l=1}^M y_{jl} (x_{[l]}^T a_j^D(w)) - \log \left( \sum_{l=1}^M \exp(x_{[l]}^T a_j^D(w)) \right) \right]$$

"logistic regression"

## Fundations of Optimization

$$f: D \subset \mathbb{R}^n \rightarrow \mathbb{R}$$

"local minimizer" "global minimizer"

"strict local minimizer"

"isolated local minimizer"

# Constrained Optimization Problem

-6-

$$\min_{x \in \Omega} f(x)$$

where  $\Omega \subset D \subset \mathbb{R}^n$  is a closed set

"local solution"

"global solution"

relation with unconstrained

$$\min f(x) + I_{\Omega}(x)$$

$$I_{\Omega}(x) = \begin{cases} 0 & \text{if } x \in \Omega \\ +\infty & \text{otherwise} \end{cases}$$

Convexity

"convex set"

$$x, y \in \Omega$$

$$\Rightarrow (1-\alpha)x + \alpha y \in \Omega \quad \forall \alpha \in [0, 1]$$

"Supporting hyperplane for  $\Omega$  at  $\bar{x} \in \Omega$ "

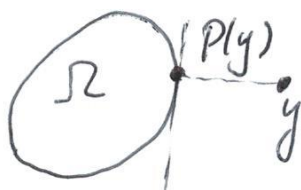
is defined by  $g \in \mathbb{R}^n$   $g \neq 0$  s.t

$$g^T (x - \bar{x}) \leq 0 \quad \text{for all } x \in \Omega$$

Projection Operator

$$P: \mathbb{R}^n \rightarrow \Omega$$

$$P(y) = \arg \min_{z \in \Omega} \|z - y\|_2^2$$



Convex function  $\phi: \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\pm\infty\}$

$$\phi((1-\alpha)x + \alpha y) \leq (1-\alpha)\phi(x) + \alpha\phi(y)$$

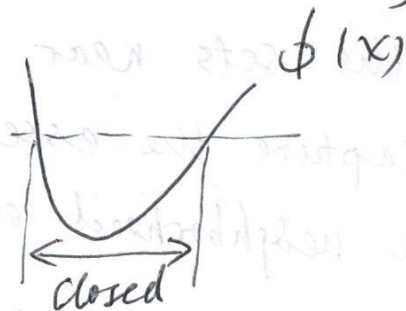
$$\forall x, y \in \mathbb{R}^n, \forall \alpha \in [0, 1]$$

"effective domain" =  $\{x \in \mathcal{D} : \phi(x) < +\infty\}$

"epigraph" =  $\text{epi } \phi := \{(x, t) \in \mathcal{D} \times \mathbb{R} : t \geq \phi(x)\}$

"proper convex function"

"closed proper convex function"

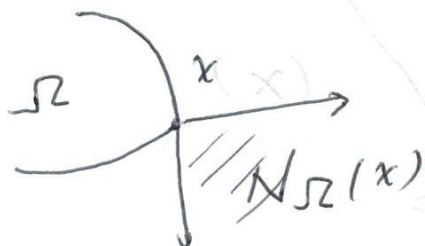


Definition "Normal Cone"

$\mathcal{D} \subset \mathbb{R}^n$  is convex set

$N_{\mathcal{D}}(x)$  = normal cone at  $\forall x \in \mathcal{D}$

$$= \{d \in \mathbb{R}^n : d^T(y-x) \leq 0 \text{ for all } y \in \mathcal{D}\}$$



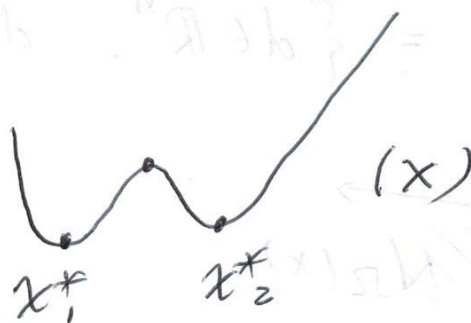
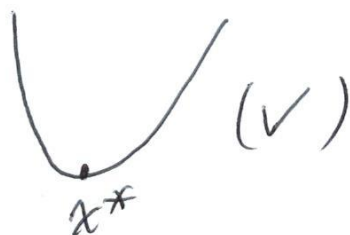
Theorem If  $\Omega_i$ ,  $i=1, 2, \dots, m$  are convex sets and  $\Omega = \bigcap_{i=1, 2, \dots, m} \Omega_i$ , then  $\forall x \in \Omega$

$$N_{\Omega}(x) \supset N_{\Omega_1}(x) + N_{\Omega_2}(x) + \dots + N_{\Omega_m}(x)$$

For " $=$ " in the above we need constraint qualifications: a linear approximation of the sets near the point in question needs to capture the essential geometry of the set itself in a neighborhood of the point.

Theorem If  $f$  is convex and  $\Omega$  closed convex then for  $\min_{x \in \Omega} f(x)$  we have

- (a) any local solution is a global solution
- (b) the set of global solutions form a convex set.





Important quantities

-9-

"modulus of continuity"  $m$  for strongly convex  $\phi$

$m > 0 \quad \forall x, y \in \text{domain of } \phi$

$$\phi((1-\alpha)x + \alpha y) \leq (1-\alpha)\phi(x) + \alpha\phi(y) - \frac{1}{2}m\alpha(1-\alpha)\|x-y\|_2^2 \quad (*)$$

Theorem (Taylor's formula)

$f: \mathbb{R}^n \rightarrow \mathbb{R}$  continuously differentiable

$x, p \in \mathbb{R}^n$

$$f(x+p) = f(x) + \int_0^1 \nabla f(x + \gamma p)^T p d\gamma$$

$$f(x+p) = f(x) + \nabla f(x + \gamma p)^T p \quad \text{for some } \gamma \in (0,1)$$

$f: \mathbb{R}^n \rightarrow \mathbb{R}$  twice continuously differentiable

$$\nabla f(x+p) = \nabla f(x) + \int_0^1 \nabla^2 f(x + \gamma p) p d\gamma$$

$$f(x+p) = f(x) + \nabla f(x)^T p + \frac{1}{2} p^T \nabla^2 f(x + \gamma p) p \quad \text{for some } \gamma \in (0,1)$$

Lipschitz constant  $L$  for  $\nabla f$

$$\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\| \quad (**)$$

for all  $x, y \in \text{dom}(f)$

### Theorem

(1) If  $f$  is continuously differentiable and convex

then 
$$f(y) \geq f(x) + (\nabla f(x))^T (y - x)$$

for  $\forall x, y \in \text{dom}(f)$

(2) If  $f$  is differentiable and strongly convex then

$$f(y) \geq f(x) + (\nabla f(x))^T (y - x) + \frac{\mu}{2} \|y - x\|^2$$

for  $\forall x, y \in \text{dom}(f)$

(3) If  $\nabla f$  is uniformly Lipschitz continuous with Lipschitz constant  $L$  and  $f$  is convex then

$$f(y) \leq f(x) + (\nabla f(x))^T (y - x) + \frac{L}{2} \|y - x\|^2$$

for  $\forall x, y \in \text{dom}(f)$

Proof. (1)  $\partial f(x) = \{ \nabla f(x) \}$

$$z \rightarrow x \quad f(z) \geq f(x) + (\nabla f(x))^T (z - x)$$

$$\wedge$$

$$\alpha f(y) + (1 - \alpha) f(x) \quad (0 < \alpha < 1)$$

$$\text{So } \alpha f(y) \geq \alpha f(x) + (\nabla f(x))^T (z - x)$$

$$\Rightarrow f(y) \geq f(x) + (\nabla f(x))^T \left( \frac{z - x}{\alpha} \right)$$

$$\parallel$$

$$(y - x)$$

(2) follows (\*)

(3) By Taylor expansion

$$f(y) - f(x) - (\nabla f(x))^T (y - x) = \int_0^1 [\nabla f(x + \gamma(y - x)) - \nabla f(x)]^T (y - x) d\gamma$$

$$\leq \int_0^1 \|\nabla f(x + \gamma(y - x)) - \nabla f(x)\| \cdot \|y - x\| d\gamma$$

$$\leq \int_0^1 L \gamma \|y - x\|^2 d\gamma = \frac{L}{2} \|y - x\|^2$$

Theorem  $f \in C^2(\mathbb{R}^n)$

$f$  is strongly convex with modulus of convexity  $m$

$$\Leftrightarrow \nabla^2 f(x) \succeq mI \quad \text{for all } x$$

$\nabla f$  is Lipschitz continuous with Lipschitz constant  $L$

$$\Leftrightarrow \nabla^2 f(x) \preceq LI \quad \text{for all } x$$

Theorem If  $f$  is differentiable and strongly convex with modulus of convexity  $m$ . Then minimizer  $x^*$  of  $f$  exists and is unique.

Key to the proof ①.  $\{x \mid f(x) \leq f(x^0)\}$  for any  $x^0$  is closed and bounded

②  $x^*$  is unique.



Theorem  $f$  is convex,  $\nabla f$  with Lipschitz constant  $L$  then  $\forall x, y \in \text{dom}(f)$

$$f(x) + (\nabla f(x))^T (y-x) + \frac{1}{2L} \|\nabla f(x) - \nabla f(y)\|^2 \leq f(y)$$

$$\frac{1}{L} \|\nabla f(x) - \nabla f(y)\|^2 \leq (\nabla f(x) - \nabla f(y))^T (x-y) \leq L\|x-y\|^2$$

If in addition  $f$  is strongly convex and with modulus of convexity  $m$ , unique minimizer  $x^*$

then 
$$f(y) - f(x) \geq -\frac{1}{2m} \|\nabla f(x)\|^2$$

$$\forall x, y \in \text{dom}(f)$$

Proof. Define  $\phi(y) = f(y) - (\nabla f(x))^T y$

$\phi$  is convex  $\nabla \phi(y) = \nabla f(y) - \nabla f(x)$

$$\nabla \phi(x) = \nabla f(x) - \nabla f(x) = 0$$

so  $x$  is a minimizer of  $\phi$

so  $\phi(x) \leq \phi(y - \frac{1}{L} \nabla \phi(y))$

$$\leq \phi(y) + (\nabla \phi(y))^T \left[-\frac{1}{L} \nabla \phi(y)\right] + \frac{L}{2} \left\| \left(-\frac{1}{L}\right) \nabla \phi(y) \right\|^2$$

$$= \phi(y) - \frac{1}{2L} \|\nabla \phi(y)\|^2$$

so  $f(x) - (\nabla f(x))^T x$

$$\leq f(y) - (\nabla f(x))^T y - \frac{1}{2L} \|\nabla f(y) - \nabla f(x)\|^2$$

i.e.  $f(y) \geq f(x) + (\nabla f(x))^T (y-x) + \frac{1}{2L} \|\nabla f(y) - \nabla f(x)\|^2$

same way  $f(x) \geq f(y) + (\nabla f(y))^T (x-y) + \frac{1}{2L} \|\nabla f(x) - \nabla f(y)\|^2$

$$\Rightarrow [(\nabla f(x))^T - (\nabla f(y))^T](x-y) \geq \frac{1}{L} \|\nabla f(x) - \nabla f(y)\|^2$$

Finally by (\*)

$$f(y) - f(x) \geq (\nabla f(x))^T (y-x) + \frac{m}{2} \|y-x\|^2$$

$$= \frac{1}{2m} \|\nabla f(x)\|^2 + (\nabla f(x))^T (y-x) + \frac{m}{2} \|y-x\|^2$$

$$= \frac{1}{2m} \|\nabla f(x)\|^2$$

$$= \frac{m}{2} \left\| y-x + \frac{1}{m} \nabla f(x) \right\|^2 - \frac{1}{2m} \|\nabla f(x)\|^2$$

$$\geq -\frac{1}{2m} \|\nabla f(x)\|^2 \quad \#$$

$$\|\nabla f(x)\|^2 \geq 2m [f(x) - f^*], m > 0$$

"generalized strong convexity condition"

"quadratic surrogate"

$$f(x) - f(x^*) = \frac{1}{2} (x - x^*)^T \nabla^2 f(x^*) (x - x^*) + o(\|x - x^*\|)$$

Optimality conditions for smooth unconstrained problem:

Theorem (Necessary Conditions for smooth unconstrained optimization)

- (a).  $f$  is continuously differentiable,  $x^*$  - local minimizer of  $\min_{x \in \mathbb{R}^n} f(x)$  then  $\nabla f(x^*) = 0$  (first-order necessary condition)
- (b).  $f$  is twice continuously differentiable,  $x^*$  - local minimizer of  $\min_{x \in \mathbb{R}^n} f(x)$  then  $\nabla f(x^*) = 0$  and  $\nabla^2 f(x^*)$  is positive semidefinite (second-order necessary condition)



When  $f$  is convex, first-order necessary condition becomes sufficient

Theorem  $f$  is continuously differentiable and convex

$$\nabla f(x^*) = 0 \Rightarrow x^* \text{ is global minimizer of } \min_{x \in \mathbb{R}^n} f(x)$$

$f$  is strongly convex  $\Rightarrow x^*$  is unique

$$\begin{aligned} \text{Key } f(y) &\geq f(x^*) + (\nabla f(x^*))^T (y - x^*) \\ &= f(x^*) \end{aligned}$$

When  $f$  is non-convex

Theorem (Second-order sufficient condition)

If  $f$  is twice continuously differentiable and that for some  $x^*$  we have  $\nabla f(x^*) = 0$  and  $\nabla^2 f(x^*)$  is positive definite

Then  $x^*$  is a strict local minimizer of

$$\min_{x \in \mathbb{R}^n} f(x).$$



Optimality conditions for smooth constrained problem  
= nonsmooth problems

$$\min_{x \in \mathbb{R}^n} [f(x) + I_{\Omega}(x)]$$

Theorem Let  $\Omega$  be closed and convex in  $\mathbb{R}^n$   
Let  $f$  be convex and differentiable

$x^*$  is a minimizer of  $\min_{x \in \mathbb{R}^n} [f(x) + I_{\Omega}(x)]$

$$\Leftrightarrow -\nabla f(x^*) \in N_{\Omega}(x^*)$$

$$\boxed{\partial I_{\Omega}(x^*) = N_{\Omega}(x^*)} \text{ (key)}$$

$$\forall d \in \partial I_{\Omega}(x^*) \quad x^* \in \Omega$$

$$I_{\Omega}(x) \geq I_{\Omega}(x^*) + d^T(x - x^*)$$

$$\text{so } d^T(x - x^*) \leq 0 \text{ of } x, x^* \in \Omega$$