

today: catalyst → accelerate  
 non-convex opt.  
 submodular opt.

Catalyst algorithm [Lin, Maillard & Hachchaoui NeurIPS 2015]

"meta-algorithm": outer loop which uses a linearly convergent alg. in inner loop to get overall acceleration (?)

main idea: use the accelerated proximal point algorithm

with approximation inner loop of prox operator

proximal point algorithm: is proximal gradient with  $f=0$

$$w_{t+1} = \text{prox}_\gamma^\Omega(w_t) \quad (\text{to solve } \min_w \Omega(w))$$

Catalyst alg. (for  $\mu$ -strongly  $F(w)$ )

let  $\rho \triangleq \frac{\mu}{\mu + \frac{1}{\gamma}}$

( $\gamma$  is algorithmic parameter)

repeat:

$$w_{t+1} \approx \underset{w}{\text{argmin}} F(w) + \frac{1}{2\gamma} \|w - z_t\|_2^2$$

to be specified

$$\text{s.t. } G_t(w_{t+1}) - \min_w G_t(w) \leq \epsilon_t$$

use inner loop optimization with warm start [eg. SAGA or AFW]

$$z_{t+1} = w_{t+1} + \beta_{t+1} (w_{t+1} - w_t)$$

"extrapolation" / "momentum"

[accelerated Nesterov trick piece]

$\beta_{t+1}$  is found using fancy equations so that everything works

• solve for  $\alpha_{t+1}$  in eq.:  $\alpha_{t+1}^2 = (1 - \alpha_{t+1}) \alpha_t^2 + \rho \alpha_{t+1}$

$$\beta_{t+1} \triangleq \frac{\alpha_t (1 - \alpha_t)}{\alpha_t^2 + \alpha_{t+1}}$$

(pick  $\alpha_{t+1} \in ]0, 1[$ )

eg.  $\alpha_0 = \sqrt{\rho}$

Catalyst trick: use  $\gamma \frac{1}{\epsilon_t}$

s.t. overall # of inner loop calls gives an overall acceleration

with clever analysis of warm starting

acceleration results:

↳ strong convexity  $G_t(w)$

acceleration results:

if unimodal alg. has a convergence  $\exp(-\frac{\tilde{\mu}}{L} t)$   $\tilde{\mu} \geq \mu + \frac{1}{\delta}$

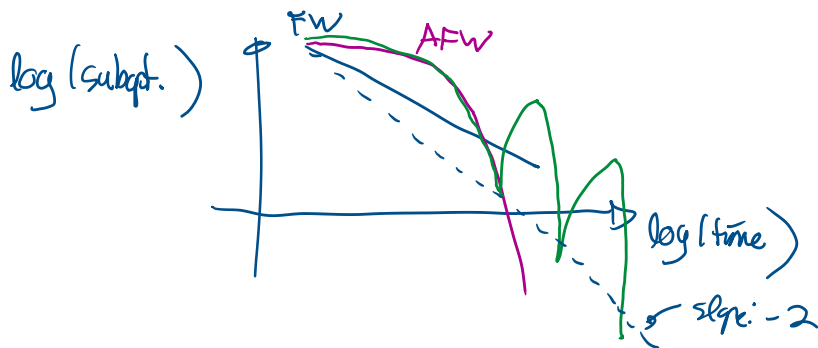
then with correct constants for  $\delta$  & schedule for  $\epsilon_t$

( $\mu$ -strongly convex  $F$ ) before: linear rate  $\rho = \frac{1}{K}$  becomes with catalyst  $\approx \frac{1}{\sqrt{K}}$  for catalyst

( $F$  convex case)  $O(\frac{1}{\epsilon})$  on  $F$  becomes  $O(\frac{1}{\epsilon^2})$

result: we can get (theory) accelerated SAGA  
 " SVRG  
 " AFW  
 etc.

Issue: catalyst is not adaptive to local strong convexity  
and slowly for choice of  $\delta, \epsilon_t, \mu$  etc...



10h15  
 10h35

non-convex optimization

recall: FW with line search on non-convex  $f$  min SST  $g(w_t) \leq O(\frac{1}{\sqrt{\epsilon}})$   
FW gap

convex:  $\mathbb{E} f(x_t) - f^* \leq \epsilon$

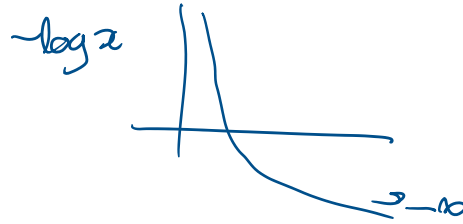
GD  $\rightarrow \frac{1}{\epsilon}$      Nesterov  $\frac{1}{\sqrt{\epsilon}}$

non-convex:  $\mathbb{E} \|Df(w_t)\|_2^2 \leq \epsilon$

a) btw: if  $f$  is  $\mu$ -strongly convex

$\Rightarrow f(w_t) - f^* \leq \frac{1}{2\mu} \|Df(w_t)\|_2^2$

b) note:  $\|\nabla f(w_t)\|$  small  $\nRightarrow f(w_t) - f^*$  is small when  $f$  is not strongly convex



(\*) can get a  $O(\frac{1}{\epsilon})$  complexity for gradient descent

$$f(w) \leq f(w_t) + \langle \nabla f(w_t), w - w_t \rangle + \frac{L}{2} \|w - w_t\|^2 \quad \forall w$$

$$w_{t+1} = w_t - \frac{1}{L} \nabla f(w_t)$$

$$\Rightarrow f(w_{t+1}) \leq f(w_t) - \frac{1}{2L} \|\nabla f(w_t)\|^2$$

(\*) key: smoothness of  $f$  but no need of  $f$  can nearby

if  $f^*$  is finite

$$\Rightarrow f(w_{t+1}) - f^* \leq f(w_0) - f^* - \frac{1}{2L} \sum_{s=0}^t \|\nabla f(w_s)\|^2$$

$$\Rightarrow \sum_{s=0}^t \|\nabla f(w_s)\|^2 \leq 2L (f(w_0) - f^*)$$

$$(t+1) \cdot \min_{s \leq t} \|\nabla f(w_s)\|^2 \leq 2L (f(w_0) - f^*)$$

$$\text{i.e. } \boxed{\min_{0 \leq s \leq t} \|\nabla f(w_s)\|^2 \leq \frac{2L}{t+1} (f(w_0) - f^*)}$$

NeurIPS 2016 tutorial "Large-Scale Optimization: Beyond Stochastic Gradient Descent and Convexity"  
[Suvri Sra slides](#)

### Faster nonconvex optimization via VR

(Reddi, Hefny, Sra, Póczos, Smola, 2016; Reddi et al., 2016)

Algorithm	Nonconvex (Lipschitz smooth)
SGD	$O(\frac{1}{\epsilon^2})$
GD	$O(\frac{n}{\epsilon})$
SVRG	$O(n + \frac{n^{2/3}}{\epsilon})$
SAGA	$O(n + \frac{n^{2/3}}{\epsilon})$
MSVRG	$O(\min(\frac{1}{\epsilon^2}, \frac{n^{2/3}}{\epsilon}))$

#### Remarks

$$\mathbb{E}[\|\nabla g(\theta_t)\|^2] \leq \epsilon$$

New results for convex case too; additional nonconvex results  
 For related results, see also (Allen-Zhu, Hazan, 2016)

## Linear rates for nonconvex problems

$$\min_{\theta \in \mathbb{R}^d} g(\theta) = \frac{1}{n} \sum_{i=1}^n f_i(\theta)$$

The **Polyak-Łojasiewicz (PL)** class of functions

$$g(\theta) - g(\theta^*) \leq \frac{1}{2\mu} \|\nabla g(\theta)\|^2$$

(Polyak, 1963); (Łojasiewicz, 1963)

## Linear rates for nonconvex problems

$$g(\theta) - g(\theta^*) \leq \frac{1}{2\mu} \|\nabla g(\theta)\|^2 \quad \Bigg| \quad \mathbb{E}[g(\theta_t) - g^*] \leq \epsilon \quad \text{😎}$$

Algorithm	Nonconvex	Nonconvex-PL
SGD	$O\left(\frac{1}{\epsilon^2}\right)$	$O\left(\frac{1}{\epsilon^2}\right)$
GD	$O\left(\frac{n}{\epsilon}\right)$	$O\left(\frac{n}{2\mu} \log \frac{1}{\epsilon}\right)$
SVRG	$O\left(n + \frac{n^{2/3}}{\epsilon}\right)$	$O\left(\left(n + \frac{n^{2/3}}{2\mu}\right) \log \frac{1}{\epsilon}\right)$
SAGA	$O\left(n + \frac{n^{2/3}}{\epsilon}\right)$	$O\left(\left(n + \frac{n^{2/3}}{2\mu}\right) \log \frac{1}{\epsilon}\right)$
MSVRG	$O\left(\min\left(\frac{1}{\epsilon^2}, \frac{n^{2/3}}{\epsilon}\right)\right)$	—

Variant of **nc-SVRG** attains this fast convergence!

(Reddi, Hefny, Sra, Póczos, Smola, 2016; Reddi et al., 2016) 22

## Submodular optimization

submodularity is an analog of convexity for tractability of set functions  
(combinatorial opt.)

$$F: 2^V \rightarrow \mathbb{R}$$

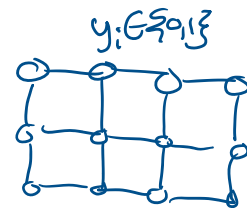
$V = \{1, \dots, d\}$  is "ground set"

$2^V = \{V \rightarrow \{0,1\}\}$  = set of all subsets of  $V$

concrete example:

Ising model

$$E(y) = \sum_i \theta_i y_i - \sum_i \sum_{j \in \text{neighbor of } i} G_{ij} y_i y_j$$



set encoding:  $A_y = \{i: y_i=1\}$   
 $F(A_y) =$

When  $G_{ij} > 0 \Rightarrow E(y)$  is submodular  
 "attractive potential"

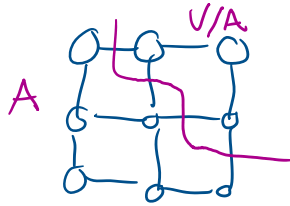
MRF here is "associative MRF"

(can minimize  $E(y)$  (or  $F(A_y)$ ) by using "graph cut algorithm")



equivalence with min cost network flow

can minimize  $f(y)$  (or  $F(A_y)$ ) by using greedy algorithm



equivalence with min cost network-flow problem

$F$  is submodular  $\Leftrightarrow F(A) + F(B) \geq F(A \cap B) + F(A \cup B) \quad \forall A, B \in \mathcal{V}$   
 $\Leftrightarrow$  function  $A \mapsto F(A \cup kE) - F(A)$  is non-increasing for all  $k$   
 "diminishing return property"

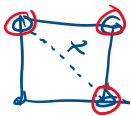
$\Rightarrow$  intuitively, that greedy alg are not "too bad" for maximization

\*  $F(A) \leq g(|A|)$  if  $g$  is concave  $\Rightarrow$  then  $F$  is submodular  
↑  
cardinality

\* link with convexity  $\rightarrow$  Lovasz extension (cf. fol.)

\* embed sets as corners of hypercube in dimension  $d \quad \forall A \in \mathcal{V} \quad v(A) = \mathbb{1}_A \in \{0,1\}^d$

Lovasz extension  $f$  extends  $F(\cdot)$  from corners to entire hypercube using convex interpolation



$f(w) = F(A_w)$  when  $w = v(A)$

let's say  $w = \sum_i \alpha_i v(A_i)$   $\Rightarrow f(w) = \sum_i \alpha_i F(A_i)$   
↑  
 $v(A_i)$

(precourse action  $f$  on  $\{0,1\}^d$ )

$F$  is submodular  $\Leftrightarrow$  Lovasz extension  $f$  is convex

(turns out)

\* can write  $f(w) = \max_{S \in B(F)} \langle s, w \rangle$   
↑  
"base polytope"

$\leftarrow$  this can be computed efficiently using greedy alg.

(LMO over  $B(F)$  is efficient)

$\min_{A \in \mathcal{V}} F(A) = \min_{w \in \{0,1\}^d} \left( \max_{S \in B(F)} \langle s, w \rangle \right)$   
↑  
 $f(w)$

$\rightarrow$  use projected subgradient method

$\mathcal{F}(w_t) = \arg \max_{S \in B(F)} \langle s, w_t \rangle$

\* with  $\ell_2$  regularization, use duality to get a smooth obj.

$\min_{S \in B(F)} \frac{1}{2} \|s\|^2$

$\rightarrow$  use "min-norm pt." alg.

you - min-norm

JEUT)

→ use "min-norm pt." alg,  
variant of FFW alg.

⊕ SOTA for submodular  
min.