

Sum of Square

For constant D, solved by interior method

$\min_{x \in D} f(x) \rightarrow \min_{y \in D'} g(y), r(y) = x$

Momentum

Smoothness:  $\| \nabla f(x) - \nabla f(y) \| \leq L \| x - y \|$   
 Adjust the lr automatically  
 Weighted sum (average of past gradient)

$f(x) = \sum_{i \in [d]} C_i \prod x_i^{\alpha_i}$

$X_{\vec{a}} = \prod_{i \in [d]} x_i^{\alpha_i}, X = \{ X_{\vec{a}} \}_{\vec{a}}$

$G(x) = \sum_{\vec{a}} C_{\vec{a}} X_{\vec{a}}$ , linear, convex

$h(x) = \sum_{\vec{a}} C(\vec{a}) \prod x_i^{\alpha_i}$

add constraint  $\sum C(\vec{a}) X_{\vec{a}} \geq 0$

> polytope constraint

Theorem:  $\min_x G(x) = \min_x f(x)$

$x^* = (x_1^*, \dots, x_d^*)$

Pseudo-expectation:  $X_{\vec{a}} \star$

$\tilde{E}_{x \sim q} = \sum_{\vec{a}} C(\vec{a}) X_{\vec{a}} \star$

Write the minimization:  $\min \tilde{E}_{x \sim q}$

for all polynomial h:  $\tilde{E}_{x \sim q} h^2 \geq 0$

Infinite h?  $\rightarrow$  only for  $\deg(h^2) \leq D$   
 $\| \vec{a} \|_1 \geq D, C(\vec{a}) \geq 0$

Let  $M_D(x)$  be the matrix, that

$[M_D(x)]_{\vec{a}, \vec{b}} = X_{\vec{a} + \vec{b}}, \| \vec{a} \|_1, \| \vec{b} \|_1 \leq D$

$M_D(x) \geq 0, \sum_{\vec{a}} C(\vec{a}) X_{\vec{a}} = \vec{1}^T M_D(x) \vec{1}$

for  $C = C(\vec{a}) \geq 0$ , degree D SOS

Constrained:  $\min f(x) \text{ s.t. } \{ h_i(x) \geq 0 \}_{i \in [m]}$

still:  $u(x) = \sum C_i X_i$

$\sum C_i h_i(x) \geq 0$

Robust Linear Regression:

$\min_{\vec{w}} \min_{\vec{x} \in \mathcal{D}} \sum d_i \langle \vec{x}^{(i)}, \vec{w} \rangle - y^{(i)}$

$d_i = d_i, \sum d_i \geq 1 - \gamma$

Interior point method

Motivation: Constraint optimization

theoretically fast

$R(x) = \begin{cases} = +\infty, & x \in D \\ \in (-\infty, \infty), & x \in D^c \end{cases}$

minimize  $f(x) + \lambda R(x)$

(can not be smooth/Lipschitz)

differentiable.

Lipschitzness of h:  $\| \nabla h(x) \|_2$

smoothness of h:  $\| \nabla^2 h(x) \|$  spectral.

Self-concordance with parameter  $\nu$

$\langle \nabla R(x), v \rangle^2 \leq \nu \langle v, \nabla^2 R(x) v \rangle$

$|\nabla^T \nabla R(x) v|^{3/2} \geq \frac{1}{2} x^T \nabla^2 R(x) (v, v)$

$\nabla R(x)(v, v) = \frac{d^2}{dt^2} R(x+tv) |_{t=0}$

$\langle \nabla R(x), v \rangle = \frac{d}{dt} R(x+tv) |_{t=0}$

$\nabla^T \nabla R(x) v = \frac{d^2}{dt^2} R(x+tv)$

change slowly  $\checkmark$  curve.

are not scaling invariant

conver  $\checkmark$ , every convex D,  $\nu$

$R_1, R_2, \nu$  self-concordant

$R_1 + R_2, \geq \nu; R_3 = R_1(Ax+b), \nu$

Example:  $D = \{ x \in \mathbb{R}^d \mid \forall i \in [d], \langle a_i, x \rangle \leq b_i \}$

$R(x) = -\sum_{i \in [d]} \log(b_i - \langle a_i, x \rangle)$

interior point method: 1.  $\min_{w \in \mathcal{S}} \sum_{i \in \mathcal{S}} (h(x_i, w) + R(w))$

$x_{t+1} = \lambda t + (1-\lambda) \arg \min_{w \in \mathcal{S}} \{ f(w) + \lambda t R(w) \}$

by pre-conditioned gd with  $\gamma$

Dikin's Ellipsoid:

$\mathcal{E}_t = \{ x \in \mathbb{R}^d \mid \langle x - x_t^*, \nabla^2 R(x_t^*) (x - x_t^*) \rangle \leq \frac{1}{4} \}$

$\frac{1}{4} \nabla^2 R(x_t^*) \leq \nabla^2 f_{t+1}(x) / \lambda_{t+1} \leq 4 \nabla^2 R(x_t^*)$

Self-concordant implies sandwich in DE

Adagrad

$x_{t+1} = x_t - M^{-1} \nabla f(x_t)$

M: diagonal pre conditioning coordinate of  $\nabla f(x_t)$  large

scale it down

each coordinate has abs value one

gradient sign:

$x_{t+1} = x_t - \gamma \text{sign}(\nabla f(x_t))$

Adagrad:  $x_{t+1} = x_t - \gamma M^{-1} \nabla f(x_t)$

$M_t = \text{diag}(\sqrt{\sum_{s=1}^t |\nabla f(x_s)|^2})$

more stable, compared to especially SGD

convergence  $\checkmark$

no need to "j" "n"

Adagrad  $\Leftrightarrow$  minor descent

$\phi_t(x) = \frac{1}{2} x^T M_t x, \nabla \phi_t(x) = M_t x$

$\nabla \phi_t(x_{t+1}) = \nabla \phi_t(x_t) - \gamma \nabla f(x_t)$

Adam: Adagrad + momentum.

No convergence even in convex

$g_{t+1} = \gamma g_t + (1-\gamma) \nabla f(x_t)$

$S_{t+1} = \beta S_t + (1-\beta) [\nabla f(x_t)]^2$

$x_{t+1} = x_t - \gamma \text{diag}(S_{t+1})^{-1/2} g_{t+1}$

Distributed

1.  $\min_{w \in \mathcal{S}} \frac{1}{n} \sum_{j \in \mathcal{S}} (h(x_j, w) + R(w))$

2.  $\nabla f(w) = \sum_{j \in \mathcal{S}} \nabla_j + \nabla R(w)$

per iter  $O(m \log)$ ,  $w \in \mathbb{R}^d$ , gd converges

ADMM: does not depend on smoothness / Lipschitzness

$\min_{w \in \mathcal{S}} \frac{1}{m} \sum_{j \in \mathcal{S}} (f_j(w_j) + \lambda \|w_j - w\|_2)$

$x_{t+1} = x_t - \frac{d}{L} \nabla f(x_t)$

scale with  $\frac{1}{L}$ ,  $L = \sum \| \nabla f_j - \nabla f \|^2$

Locally:  $w_j = \arg \min (f_j(w_j) + \lambda \|w_j - w\|_2 + c_j)$

$w^{t+1} = \frac{\sum_j d_j^{(t)}}{2m\lambda} + \frac{1}{\lambda} \sum w_j^{(t+1)}$

$d_j^{t+1} = d_j^t - \gamma (w_j^{(t+1)} - w_j^{(t)})$   $O(m \log)$

Online optimization

env: pick  $f_t = D \rightarrow R$

act: choose  $x_t \in D$ , not knowing  $f_t$

env:  $f_t(x_t), \nabla f_t(x_t)$

env changing independent of act.

want act to adapt

$\sum_{t \in [T]} f_t(x_t) \min_{x \in D} \sum_{t \in [T]} f_t(x)$

do as good as best fixed x

regret of player:

$R = \frac{1}{T} \sum_{t \in [T]} f_t(x_t) - \min_{x \in D} \sum_{t \in [T]} f_t(x)$

if  $f_t = f$ , gradient descent, convex

if  $E[f_t(x)] = f(x)$  SGD

when  $f_t$  are convex: gd  $\checkmark$

L-Lipschitz

$x^* = \arg \min \sum f_t(x), \gamma = \frac{\|x^* - x_0\|}{L T}$

$R = \frac{1}{T} \sum f_t(x_t) - \frac{1}{T} \sum f_t(x^*) \leq \frac{\|x^* - x_0\|_2}{\sqrt{T}}$

projected:  $x_{t+1} = \Pi_D (x_t - \gamma \nabla f_t(x_t))$

three-term MD

$f_t(x_t) \leq f_t(x) + \frac{1}{2} \gamma^2 (\|x - x_t\|_2^2 - \|x - x_{t+1}\|_2^2)$

$\frac{1}{T} \sum f_t(x_t) \leq \frac{1}{T} \sum f_t(x) + \frac{1}{2T} \sum \|x - x_t\|_2^2$

Stokes: for every  $r > 0, \nu$

randomly sampled unit vector

$E[\nu^T f(x_t + \nu)] = \nu^T \nabla f(x)$

$x_{t+1} = x_t - \frac{d}{L} \nabla f(x_t)$

scale with  $\frac{1}{L}$ ,  $L = \sum \| \nabla f_j - \nabla f \|^2$

Simulated Annealing

$f: \mathbb{R}^d \rightarrow \mathbb{R}$   
 $p(x) \propto e^{-\lambda f(x)}$   
 $\lambda \gg 0$ , uniform temperature  $\rightarrow \infty$   
 $\lambda \gg \infty$ , global minimizer  $f$

Metropolis-Hastings

$x_{t+1} \sim N(x_t, \sigma^2 I)$   
 $\alpha = \min \left\{ \frac{p(y_{t+1})}{p(x_t)}, 1 \right\}$   
 $x_{t+1} = \begin{cases} y_{t+1}, & \text{w.p. } \alpha \\ x_t, & \text{w.p. } 1-\alpha \end{cases}$

theorem:  $t \rightarrow \infty, x_t \rightarrow \sim p(x)$

Stationary distribution of Metropolis-Has

$$\begin{aligned} \pi(x) &= \int_{x'} p(x|x') \pi(x') dx' \\ &= \int_{x'} \alpha_{xx'} n(x|x') \pi(x') dx' \\ &+ \pi(x) \int_z (1-\alpha_z) n(z|x) dz \\ &= \int_{x'} n(x|x') \alpha_{xx'} \pi(x') dx' + \dots \\ &= \pi(x) \int_z n(z|x) dz = \pi(x) \end{aligned}$$

key idea:  $Q \rightarrow P$ . converge faster.

Simulate annealing

First: sample  $x_0 \sim D$ .  $\lambda_0 > 0$  small.  
 step:  $\lambda_{t+1} = (1+\eta)\lambda_t$   
 M-H  
 $p_{\lambda_{t+1}}(x) \propto e^{-\lambda_{t+1} f(x)}$

Evolve the  $p$ : random  $\rightarrow$  supports only on minimizer  
 no efficient rate guarantee, unless  $f$  convex  
 convergence  $\checkmark$ . step  $\rightarrow \infty$ . MH.

Evolutionary strategies

$f: \mathbb{R}^d \rightarrow \mathbb{R}$   
 $G_1, \dots, G_N \sim N(0, I_{d \times d})$   
 $f_i = f(x_t + \sigma_t \epsilon_i)$   
 $x_{t+1} = x_t + \frac{1}{\sqrt{\sigma_t}} \sum f_i \epsilon_i$   
 $f_i \uparrow$  higher weight for  $\epsilon_i$   
 do sGD on  $F_t = f * g_t$   
 $g_t = N(0, \sigma_t^2 I)$   
 $*: [f * g](x) = \int y f(y) g(y-x) dy$   
 $\nabla F_t = f * \nabla g_t$  stoke's.  
 escape local minima when  $\sigma_t \uparrow$   
Non-convex

Stationary distribution of Metropolis-Has

$$\begin{aligned} f(x+z) &= f(x) + \langle \nabla f(x), z \rangle + \frac{1}{2} z^T \nabla^2 f(x) z + o(\|z\|^2) \\ \nabla^2 f(x) \text{ PSD?} \\ \text{second order local minima: } \nabla f(x) = 0 \\ f(x) = x^2, x=0 \text{ local minimum } \nabla^2 f(x) > 0 \\ \text{Hessian descent, not in practice} \\ y_{t+1} &= x_t - \eta \nabla f(x_t) \end{aligned}$$

unit vector  $v$ , corresponding to eigen of  $\nabla^2 f(x)$  with smallest eigen value.

$z_{t+1,1} = x_t - \eta v$ ,  $z_{t+1,2} = x_t + \eta v$   
 $x_{t+1} = \arg \min \{ f(y_{t+1}), f(z_{t+1,1}), f(z_{t+1,2}) \}$   
 $\gamma$ -Hessian Lipschitz:

$$\begin{aligned} f(x+z) &\leq f(x) + \langle \nabla f(x), z \rangle + \frac{1}{2} z^T \nabla^2 f(x) z + \gamma \|z\|^3 \\ \frac{1}{2} (f(z_{t+1,1}) + f(z_{t+1,2})) &\leq f(x_t) - \frac{\eta^2}{4} \delta \\ \text{Hessian negative, locally high non-convex} \\ f(y_{t+1}) &\leq f(x_t) - \frac{\eta}{2} \|\nabla f(x_t)\|^2 \end{aligned}$$

Noisy GD  $\checkmark$  efficiently find

Second order minima  
 $x_{t+1} = x_t - \eta \nabla f(x_t)$   
 $x_{t+1} \leftarrow x_{t+1} + \xi_{t+1} \sim N(0, \sigma^2 I)$   
 $\nabla f(x_{t+1}) = \nabla f(x_t) - \eta \nabla^2 f(x_t) \nabla f(x_t) + O(\eta^2)$   
 $\sigma_t^2(x) = K(x, x) - \frac{V(x)^T \nabla f(x)}{\|\nabla f(x)\|^2}$   
 $\approx (I - \eta \nabla^2 f(x_t)) \nabla f(x_t)$   
 when  $\nabla f(x_t)$  not PSD.  
 $\|I - \eta \nabla^2 f(x_t)\|_2 > 1$ , gradient large.  
 Only stop when gradient small  
 Hessian PSD

$\|\nabla f(x_t)\|_2 \leq \rho \cdot \langle v, \nabla f(x_{t+1}) \rangle$  increasing acquisition function  
 geometry rate  
 convergence rate poly( $1/\epsilon, 1/\delta$ )  
 dimension free. gradient  $\leq \epsilon$   
 Hessian  $\geq \delta I$

BO

"local minima can not efficiently  
 spirit of BO: Graduate Student Descent  
 No need of computing gradient, global optimization.  
 low dimension  $\checkmark$  high inefficient.  
 vector analog of Ellipsoid algorithm.  
 Generic routine:

$f: \mathbb{R}^d \rightarrow \mathbb{R}$ , convex or not  
 step  $t$ : compute a distribution  $P_t$  over  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ . using  $f(x_0), \dots, f(x_t)$   
 find  $x_{t+1}$ , maximize acquisition func associated with  $P_t$ .

GP: given kernel  $k(x, y)$   
 $k(x, y) = \langle \alpha(x), \alpha(y) \rangle = e^{-\frac{\|x-y\|_2}{2\ell}}$   
 measures inverse distance.  
 Gaussian kernel, closer, larger

Refine  $V_t(x) = (k(x, x_0), \dots, k(x, x_t))^T$

$F_t = (f(x_0), f(x_1), \dots, f(x_t))^T$   
 $[M_t]_{i,j} = k(x_i, x_j)$   
 $P_t(x) \sim N(\mu_t(x), \sigma_t^2(x))$   
 where  $\mu_t(x) = V_t(x)^T M_t^{-1} F_t$   
 $M_t(x) = \begin{pmatrix} k(x, x) & V_t(x)^T \\ V_t(x) & M_t \end{pmatrix} \rightarrow \text{PSD}$   
 $\mu_t(x_j) = f(x_j), \sigma_t^2(x_j) = 0$   
 $x_{t+1} = \arg \max_x E_{g \sim P_t} [g(x) - f_t^*]^2$

$E_t(x) = E_{g \sim P_t} [g(x) - f_t^*]^2$   
 $[Z]^* = \text{ReLU}(Z)$   
 using (noisy) gradient descent compute maximizer cheap for  $E_t(x)$   
 optimistic exploration theory unknown

Over-parameterization  
 matrix sensing  
 $\min_{U, V} \frac{1}{N} \sum_i \langle A_i, UV^T - Y_i \rangle^2 + \lambda (\|U\|_F^2 + \|V\|_F^2)$   
 $\min_M \frac{1}{N} \sum_i \langle A_i, M - Y_i \rangle^2 + \lambda \|M\|_*$   
 non convex matrix product equivalent to nuclear norm

$\min_{W \in \mathbb{R}^d} f(W) \rightarrow \min_{W \in \mathbb{R}^d} f(W)$ ,  $D \gg d$   
 theorem. all the second order minima  $U, V$  has objective  $\leq 1 - \epsilon$ , when  $K \gg R$   
 $U, V \in \mathbb{R}^d \times \mathbb{R}^R$   
 $U^*, V^* \in \mathbb{R}^d \times \mathbb{R}^R, \xi$   
 $R \geq \text{poly}(\|Z\|_F, \alpha(\epsilon) (\|Z\|_F^2) \leq \alpha(\epsilon)$

proof:  $f(U, V) \geq 1 - \epsilon, f(U^*, V^*) \leq \epsilon$   
 $\mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$   
 $\|U - U^*\|_2 \leq \frac{\epsilon}{\alpha(\epsilon)}$   
 $(1 - \eta^2) \|U - U^*\|_2 \leq f(U, V) - \epsilon \leq \alpha(\epsilon)$   
 contradict with  $\nabla f(U, V) = 0, \nabla f(U, V) \gg 0 \rightarrow f(U, V) - \alpha(\epsilon)$

Sampling, generative model.

generate from same distribut  
 diffusion approach:  $q_0 \rightarrow q^*$   
 Brownian motion.  $X_0 \sim q^*$   
 $dX_t = -X_t dt + \sqrt{2t} dB_t \rightarrow$  Gaussian noise  
 As  $t \rightarrow \infty, X_t \rightarrow N(0, 1)$   
 $X_t = e^{-t} X_0 + \sqrt{1 - e^{-2t}}$

Forward:  $q_0 = q^*, q_T \rightarrow$  Gaussian  
 $P_t(x) = q_{T-t}(x)$  sufficient (large  $T$ )  
 satisfies,  $\frac{\partial q_t(x)}{\partial t} = \langle \nabla, X q_t(x) \rangle + \langle \nabla^2 q_t(x), x \rangle$   
 Backward satisfies  
 $\frac{\partial P_t(x)}{\partial t} = \frac{\partial q_{T-t}(x)}{\partial t} = -\langle \nabla, X q_{T-t}(x) \rangle - \langle \nabla^2 q_{T-t}(x), x \rangle$   
 due to  $\nabla \ln p_t(x) = \frac{\nabla p_t(x)}{p_t(x)}$

only  $x \sim q^* = q_0$ , score estimation to compute  $\nabla \ln q_t(x_t)$ :  
 minimize:  $\min_x E_x \|s(x_t) - \nabla \ln q_t(x_t)\|_2^2$   
 Here:  $X_0 \sim q^*, X_t = e^{-t} X_0 + \sqrt{1 - e^{-2t}}$

Stein's formula.  $\int \nabla \cdot v dy = \int \langle y, v(y) \rangle dy$   
 $y \sim N(\mu, \sigma^2)$   
 $\min_x E_x \|s(x_t) + \frac{1}{\sqrt{1 - e^{-2t}}} y\|_2^2$   
 where  $X_0 \sim q^*$  target distribution  
 $X_t = e^{-t} X_0 + \sqrt{1 - e^{-2t}}$

DPM: then use backward diffusion, starting  $Y_0 \sim N(0, 1)$   
 $dY_t = Y_t + 2S_t - t(Y_t) dt$

$E$