

Regret: and SGD

①

• at each step, we see a new objective function

• online optimization

at time step t , alg sees

$f_t: K \rightarrow \mathbb{R}$ and proposes

x_t

• how to score this?

$$\text{regret} = \sum_t f_t(x_t) - \min_{x^* \in K} \sum_t f_t(x^*)$$

time average of cost functions at chosen points.

you see all ~~of~~ f_t 's in advance - find smallest average.

The reason this is helpful is we can prove bounds, which lead to new alg.s (2)

take g.d. and substitute

$$x_{t+1} \leftarrow x_t - \eta \nabla f_t(x_t)$$

↑
use convex function
~~not~~

from proof of G-2, we have

$$f_t(x_t) + \phi_{t+1} - \phi_t \leq f_t(x^*) + \frac{1}{2} \eta G^2$$

(cause true for any convex).

Sum this

$$\begin{aligned} \sum_{t=1}^T (f_t(x_t) - f_t(x^*)) &\leq \sum_{t=1}^T (\phi_t - \phi_{t+1}) + \frac{1}{2} \eta G^2 \\ &\leq \phi_1 + \frac{1}{2} \eta T G^2 \end{aligned}$$

now plug in $T \geq \frac{\|x, -x^*\|^2 G^2}{\varepsilon^2}$

3

$\eta = \frac{\|x, -x^*\|}{G\sqrt{T}}$, do algebra. to

get.

$$\frac{1}{T} \sum_{t=0}^T (f_t(x_t) - f_t(x^*)) \leq \frac{\|x, -x^*\| G}{\sqrt{T}} \leq \varepsilon$$

↑
regret.

But pure sgd is not particularly ⁽⁴⁾
well adapted to some cases

• for regret bound, we assumed
nothing about $f_t(x)$. (except convexity)

• Now imagine we know something.

• e.g.

$$f_t(x) = \frac{1}{2} [a_t^T x]^2$$

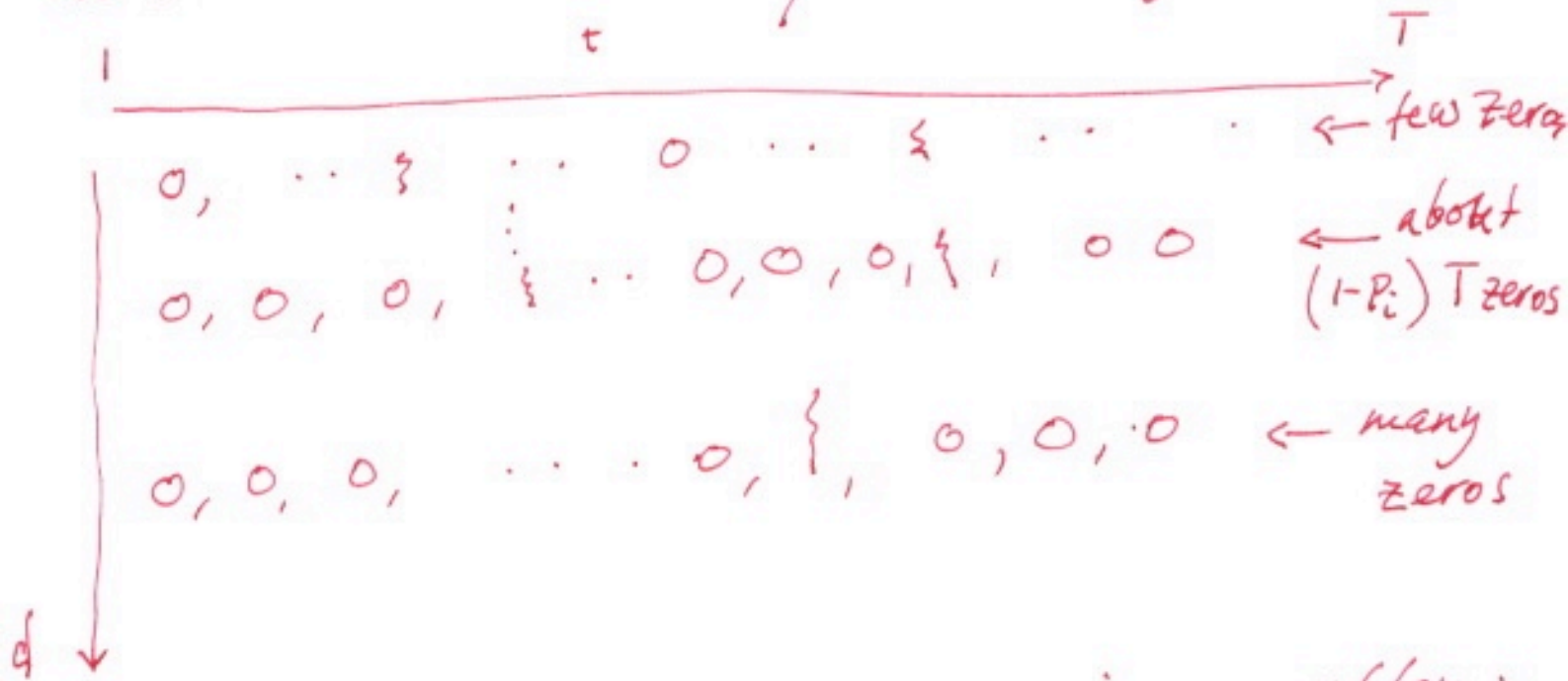
where a is very
sparse

i 'th component of a_t
has prob of non-zero
 P_i iid

and (aren't we lucky!)
 P is sorted

$$P_1 \geq P_2 \dots \geq P_d$$

think about the sequence of gradients $\{g_t\}_{t=1}^T$



In this case, we face a serious problem:

- for some components, we take a large step very occasionally
- connecting this step would be hard (slow) cause gradient would be 0 at most t
- we would like to reduce the step using some estimate of sparsity.

Example

$$\underline{x} = \begin{bmatrix} x_1 \\ \vdots \\ \vdots \\ \vdots \\ x_d \end{bmatrix}$$

$$f_t(\underline{x}) = \sum_i \delta_{it} \frac{x_i^2}{2}$$

where $\delta_{i,t} = \begin{cases} 1 & \text{with prob } p_i \\ 0 & (1-p_i) \end{cases}$

and IID over t

Now think about component i ?

iteration

$$x_i^{(\omega+1)} = x_i^{(\omega)} - \eta g_i^{(\omega)}$$

$$g_i^{(\omega)} = \begin{cases} x_i^{(\omega)} & \text{with prob } p_i \\ 0 & \text{otherwise} \end{cases}$$

so $x^{(\omega+1)} \approx (1-\eta)^{p(\omega+1)} x_0$

essentially, you see the gradient about $p(\omega+1)$ times

but this converges rather slowly (7)

IDEA :

$$x_i^{(\omega+1)} = (\text{shrink}) \left[x_i^{(\omega)} - \frac{\eta}{\text{shrink}} g_i^{(\omega)} \right]$$

where shrink uses the frequency with which $g_i^{(\omega)} \neq 0$

~~ADAGRAD~~ ADAGRAD :

Define $G_{\omega} = \sum_{t=1}^{\omega} g^{(t)} g^{(t)T}$

outer product matrix of past gradients

• Consider $\text{diag } G_{\omega} = \Delta_{\omega}$

• i'th diagonal component \bullet is approx $E[g_i^2]$

update

$$\underline{x}^{(\omega+1)} = \begin{bmatrix} \Delta^{(\omega)} \end{bmatrix}^{1/2} \left[\underline{x}^{(\omega)} - \eta \begin{bmatrix} \Delta^{(\omega)} \\ g^{(\omega)} \end{bmatrix}^{(-1/2)} \right] \quad (8)$$

notice you could have a divide by zero issue here — actually use

$$\begin{bmatrix} \Delta^{(\omega)} + \epsilon I \end{bmatrix}^{(-1/2)}$$

Fact: (Duchi et al, p2125)

let $x^{(\omega)}$ be generated as above and assume

$$\max_t \|x^* - x_t\|_\infty \leq D_\infty$$

use $\eta = \frac{D_\infty}{\sqrt{2}}$ then for x^* we

get.

9

$$R(T) \leq \sqrt{2} D_{\infty} \sum_{i=1}^d \|g_{1:T, i}\|_2$$

↑
regret at T

↑
length of i 'th
component of gradient
sequence (Note: NOT
squared length)

this is better than previous bound when

- D_{∞} is reasonable (eg, we're in a box)
- gradients have nice sparsity properties

More examples: Duchi et al 2015...

Alternative view

(10)

- $f_t(x)$ is a "noisy" function rather than an arbitrary seq of convex functions (avoid defining, as do authors!)
- Consider steps.
 - in each component, we would like the step to be bounded by about the same amount

• consider

$$\alpha \left[\frac{\text{mean}(\text{gradient})}{\sqrt{\text{mean}(\text{gradient}^2)}} \right]$$

↳ this will tend to be about 1 or -1

- this "filters" the gradient ①

ISSUE

- estimate means?
- these estimates should be
 $\text{mean}(\text{gradient}_t)$
over
all possible
 f_t
- But we can't compute that.
- instead, compute moving average

$$v_0 = 0$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

↑
element wise

- this incorporates past gradients
into est of $\text{mean}(\text{gradient}^2)$

BUT as an est of $\text{mean}(g_t^2)$ (12)
it is biased (~~Welling et al~~, p3)
Kajima + Welling

use

$$\frac{v_t}{(1 - \beta_2^t)}$$

(same arg for $\text{mean}(g_t)$)

we now have ADAM

Experience shows

- fast convergence to slightly worse solns
- SGD with good step lengths gets better solns but much slower

WHY ?

Algorithm 2: *AdaMax*, a variant of Adam based on the infinity norm. See section 7.1 for details. Good default settings for the tested machine learning problems are $\alpha = 0.002$, $\beta_1 = 0.9$ and $\beta_2 = 0.999$. With β_1^t we denote β_1 to the power t . Here, $(\alpha/(1 - \beta_1^t))$ is the learning rate with the bias-correction term for the first moment. All operations on vectors are element-wise.

Require: α : Stepsize

Require: $\beta_1, \beta_2 \in [0, 1)$: Exponential decay rates

Require: $f(\theta)$: Stochastic objective function with parameters θ

Require: θ_0 : Initial parameter vector

$m_0 \leftarrow 0$ (Initialize 1st moment vector)

$u_0 \leftarrow 0$ (Initialize the exponentially weighted infinity norm)

$t \leftarrow 0$ (Initialize timestep)

while θ_t not converged **do**

$t \leftarrow t + 1$

$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$ (Get gradients w.r.t. stochastic objective at timestep t)

$m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$ (Update biased first moment estimate)

$u_t \leftarrow \max(\beta_2 \cdot u_{t-1}, |g_t|)$ (Update the exponentially weighted infinity norm)

$\theta_t \leftarrow \theta_{t-1} - (\alpha/(1 - \beta_1^t)) \cdot m_t/u_t$ (Update parameters)

end while

return θ_t (Resulting parameters)

Regret