

Striking behavior of SGD.

①

consider

$$\min_{X \succeq 0} F(X) = \|A(X) - y\|_2^2$$

here X is a matrix, A a linear fn that accepts a matrix and makes a vector, y a vector, A : symmetric

$$A_i(X) = \text{i-th component of } A(X)$$

$$= \text{trace}(A^T X)$$

$$= \text{sum}(\text{sum}(A .* X))$$

X is $n \times n$, and we are interested in

the case where $\dim y = m \ll n^2$

(so you *really* can't recover X without something interesting going on).

We apply some opt procedure to

$F(x)$ to get x^*

1) (obvious) there are many x^* that are optimal

2) (less so) the x^* we get depends on the procedure we use

3) we are interested in

$\|x - x^*\|^2 \leftarrow \text{generalization error}$

It should be clear that this could be pretty much arbitrary (eg $A_i = \text{matrix with one's in it}$)

rather than work with X ,
consider working with u

$$X = uu^T \quad (\text{so p.s.d. is satisfied})$$

now if u is tall + thin, we are imposing a rank constraint on X

- doing so has a long history
- it is often useful

(recommender systems, feature seln, etc)

instead, for us u is square.

$$\min_u f(u) = \|A(uu^T) - y\|_2^2$$

1) Gradient descent does odd things on this problem (humarsekar, fig 1)

Why:

very rough analysis.

consider $u^{(n+1)} = u^{(n)} - \eta \nabla_u f(u)$.

- this is a form of numerical integration of a dynamical system

$$u_{(t+\Delta t)} = u_{(t)} - \Delta t \nabla_u f(u)$$

or $\dot{u}_t = \frac{du}{dt} = -\nabla_u f(u)$.

take the gradient to get.

$$\nabla_u f(u) = 4 \left[\sum_i r_i A_i \right] u$$

where

$$r_i = A_i (u u^T) - y_i$$

(5)

now

$$X_t = u_t u_t^T$$

so

$$\dot{X}_t = \dot{u}_t u_t^T + u_t \dot{u}_t^T$$

$$= -\left(\sum_i r_i A_i\right) X_t - X_t \left(\sum_i r_i A_i\right)$$

now we care about limit points of this flow.

→ assume $m=1$, $A_i = A$, $y_i = y$.

$$\dot{X} = -r_t (AX_t + X_t A)$$

$$X_t = \exp\left(\int_0^t r_t A dt\right) X_0 \exp\left(\int_0^t r_t A dt\right)$$

↑ initial cond

$$s_t = -\int_0^t r_t dt$$

Now assume it's OK to pass to (6)
DS limit

IF X_t has a limit point,
we expect $|S_t|$ to be big. (G claim ∞ ?)

Now

$$\exp(sA) = I + sA + \frac{s^2 A^2}{2} + \dots$$

if v is eigenvector, λ eval

$$\exp(sA)v = \begin{pmatrix} e^{s\lambda} \\ \cdot \end{pmatrix} v.$$

→ this means that, as long as you
start away from zero,
one ev. of X should be a lot
larger than others

gradient descent for a linear predictor

we have $\{x_n, y_n\}_{1:N}$, x_i a vector
 y_i a label $\in \{1, -1\}$.

$$\min_w \sum_{i=1}^N \ell(y_i w^T x_i)$$

is our learning problem, for ℓ some loss.

Assume: - linear separability
 $w_*^T x_i \geq 0$ all i

log-loss, NOT hinge loss

- ℓ is: positive, differentiable
monotonically decreasing to zero
derivative is Lipschitz?
 $\lim_{u \rightarrow -\infty} \ell(u) \neq 0$ (so you can go downhill)

- now assume WLOG all labels are 1 (by flipping sign of x_u as required) ⑧

- So we care about

$$\sum_i l(w^T x_i)$$

- GD gives

$$\begin{aligned} w(t+1) &= w(t) - \eta \nabla L \\ &= w(t) - \eta \sum_i l' \cdot x_i \end{aligned}$$

- notice that

$$\lim_{t \rightarrow \infty} \|w(t)\| = \infty$$

~~Q.E.D.~~

because

$$w_*^T \nabla L = \sum_i l'(w_*^T x_i)$$

but $w_*^T x_i > 0$ (lin sep)

and $l' < 0$

so $w_*^T \nabla L < 0$ so $\nabla L \neq 0$

BUT GD on smooth loss goes to critical pt.

So $\nabla L \xrightarrow{t \rightarrow \infty} 0$

which means enough t

(~~or else~~ we must have $w(t)^T x_i > 0$ for large $l' \rightarrow 0$)

So what is

$$\lim_{t \rightarrow \infty} \frac{w(t)}{\|w(t)\|} \quad ?$$

(if it exists).

work with $l(u) = e^{-u}$

assume $\frac{w(t)}{\|w(t)\|} \rightarrow w_\infty$

$$w(t) = g(t) w_\infty + r(t).$$

where $g(t) \rightarrow \infty$, $\frac{r(t)}{g(t)} \rightarrow 0$

then

$$\begin{aligned} -\nabla L(w) &= \sum_i \exp(-w(t) x_i) x_i \\ &= \sum_i e^{-g w_\infty^T x_i} e^{-r^T x_i} x_i \end{aligned}$$

Now notice that this means, (11)
 for big t , only samples w /
Smallest value of $w_{\infty}^T x_i$ contribute
 to gradient.

so w_{∞} is dominated by these

$\Rightarrow w_{\infty} =$ linear combination of x_i st $w_{\infty}^T x_i$ is min.

now consider $\hat{w} = \frac{w_{\infty}}{\min_k (w_{\infty}^T x_i)}$

support vectors
more than 1

$$\text{so } \hat{w} = \sum_i \alpha_i x_i$$

and either $\alpha_i > 0$, $\hat{w}^T x_i = 1$
 OR $\alpha_i = 0$, $\hat{w}^T x_i > 1$

(KKT for SVM!)

In each case, choice of optim procedure was not innocent from persp of learning. (12)

Q: What about ADAM, Adagrad, Momentum, etc?

Q: What about SGD?

Q: What about convnets?