

Optimization

why does it work ?

minima ? How many ?

Do they control non-complexity ?

Plan

Background on SGD and NN
- cross-entropy, traj

Minima

- Number, Berout, square loss
- Degeneracy Berout
- SGD and Langerin
- SGD finds global minima

Stability

- equil. points ?
- Hessian ?
- Variations

Loss functions

L² loss $L(w) = \sum_i^N (f(x_i; w) - y_i)^2$

Exp. loss $L(w) = \sum_i^N e^{-y_i f(x_i; w)}$

logistic is $\sum_i^N \ln(1 + e^{-y_i f(x_i; w)})$

crossentropy is multi-label version logistic

Gradient Descent optimization

$$\dot{W}_k^{ij} = -\gamma \frac{\partial}{\partial W_k^{ij}} L$$



$$W_k^{ij, t+1} - W_k^{ij, t} = -\gamma \frac{\partial L}{\partial W_k^{ij}}$$

Instead of $L = \sum_i^N l_i$

use minibatches selected at random for

each iteration \Rightarrow SGD

Minima

Can we say how many
which kind, independent of
GD?

Key fact: DNN are usually
overparametrized with $M \gg N$

$$M \text{ is } \neq W_H^{i,T}$$

$$N \text{ is } \neq \text{eqts.}$$

Bezout theorem

Instead of
$$\min_w \sum_i^N \left(f(x_i; w) - y_i \right)^2$$

consider

$$f(x_i; w_k^{(i)}) = y_i$$

$$i = 1, \dots, N \quad i, j, n \dots$$

- because it is easy to find zero error
- because overparameterization

Berout Theorem

If $\delta(z)$ were $\mathcal{P}(z)$ then
 $f(\underline{x}; w)$ is multivariate polynomial.
in the $w_k^{i,T}$ (and in the \underline{x}).

Then

$$f(x_i) - y_i = 0 \quad i = 1, \dots, N$$

is a system of N polynomial equations
in M variables. $N \sim 60k$ in CIFAR, $M \sim 300k$

Berout theorem

A set of N polynomial eqts.
in M variables of degree \mathcal{H} has
 $\leq K^N$ isolated solutions if $M = \mathcal{H} + 1$.
 $M \gg N$ then The solutions

are degenerate.

Remarks

- ① This is similar to systems of linear equations.
- ② For the size of N today \neq isolated solutions is very high \approx protos universe and furthermore degenerate because of $M \gg N$.
- ③ The point on M and N and degeneracy is what we use next.

Because $M > N$ the global minima corresponding to zero error for all x_i ($f(x_i; w) - y_i = 0 \quad i = 1, \dots, N$) are degenerate.

What about all minima?

The stationary points of the gradient are

$$-\nabla_w L = 0 \quad \text{which means}$$

$$-\frac{\partial}{\partial w_{ij}^k} \sum_i^N (f(x_i) - y_i)^2 = 0$$

These are M equations in M unknowns.

If f polynomial the equations are polynomial equations, Bezout theorem applies, the solutions are

in general not degenerate.

This \Rightarrow

① global solutions are degenerate

② local minima are not degenerate

SGD, SGDL

For the next step I need to establish similarity between SGD and Langevin equations.

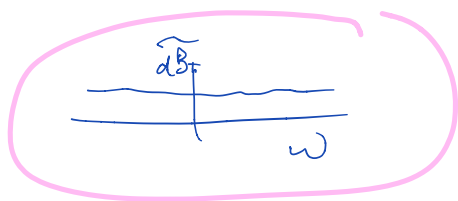
$$L = \sum_i^N (f(x_i) - y_i)^2 = \sum_i V(f, z_i)$$

$$\text{- GD } \Delta w_t = w_{t+1} - w_t = -\eta_t \nabla_w L \sim \dot{w}_n^{is} = -\frac{\partial}{\partial w_n^{is}} L$$

dynamical gradient system

$$\text{- SGD } \Delta w_t = -\eta_t \nabla_w V(f, z) \text{ with } z_i \text{ chosen at random}$$

- GDL



$$\dot{W}_k^{(0)} = - \frac{\partial L}{\partial W_k^{(1)}} + dB_t$$

SDE

where dB_t is the derivative of the Brownian motion, that is zero mean white noise with Gaussian statistics

SGD is similar to GDL in simulations and also if I write SGD as

$$\dot{W} = - \nabla_w (L - L + V) = - \nabla_w L + \xi_+$$

where $\xi_+ = \nabla_w (L - V)$ is a pseudo noise

s.t. $\mathbb{E}(\xi_+) = 0$. ξ_+ is defined in

terms of minibatches (where CLT

applies, giving ξ_+ some Gauss like

properties)

Let us speak about GDL which is a SDE.

$$\dot{w}^0 = -\nabla_w L(w) + dB_t$$

Its "solution" for stationary prob. distrib

is

$$p(w) = \frac{1}{Z} e^{-\frac{L}{T}}$$

This means that if



Important: P shows concentration of probabilities with large d ; most of probability mass is in large volume minima of Z : see slides

The conclusion is that the prob solution of GDZ prefers with high probability degenerate minima. Together with Besout conclusions this implies that GDZ prefers global minima (vs local ones). Because of GDZ \approx SGD this is valid also for SGD.

The last point in this class - which is also a harbinger of next class - is about the structure of the solutions of GD with square loss in the overparametrized ~~case~~.

The dynamical system is

$$\dot{W}_{k^i}^{i^j} = -\nabla_{W^i} L = 2 \sum_i^N \underbrace{(y_i - f(x_i))}_{E_i} \frac{\partial f(x_i)}{\partial W_{k^i}^{i^j}}$$

if $E_i = 0 \forall i$ then

$$\dot{W}_{k^i}^{i^j} = 0$$

may be zero too!

Are these solutions stable? Unique?

Let us look at Hessian of L

$$\frac{\partial^2 L}{\partial W_{k^i}^{i^j} \partial W_{k^i}^{a^b}} = 2 \sum_i^N \left\{ \frac{\partial f}{\partial W_{k^i}^{a^b}} \frac{\partial f}{\partial W_{k^i}^{i^j}} + (y_i - f(x_i)) \frac{\partial^2 f}{\partial W_{k^i}^{i^j} \partial W_{k^i}^{a^b}} \right\}$$

$$= \text{if } \bar{E} = 0 \quad \text{to } \propto \sum_i^N - \frac{\partial f}{\partial w_n^{ab}} \frac{\partial f}{\partial w_n^{ij}}$$

if $-H$ is p.d then stability. But

$\frac{\partial f}{\partial w_n^{ab}} \frac{\partial f}{\partial w_n^{ij}}$ is often \emptyset ... degenerate directions
 valleys ... as expected from Bethe analysis
 ...