

Optimization

⑤ Recap

$$1) R(f(x_i; w)) = \sum_{i=1}^N l(f(x_i; w), y_i)$$

---, one term for each example you want to label

2) in practice stochastic gradient descent / mini-batch gradient descent, K random examples

↓
(K should be not too large bc randomness helps escape spurious local minima)

reduce variance of unbiased estimator

$$3) \text{Nesterov acc } g_{t+1} = x_t - \eta \nabla f(x_t)$$

$$x_{t+1} = g_{t+1} + \mu \underbrace{(g_{t+1} - g_t)}_{\text{momentum}}$$

iterations to reduce distance to opt. by const. factor

$$\text{GD} \sim \frac{P}{\lambda}$$

$$\text{AGD} \sim \sqrt{P/\lambda}$$

$\lambda, P = \text{smoothness, strong convexity}$

4) gradient flow: continuous time limit of GD

$$\dot{x}(t) = -\nabla f(x(t))$$

5) Moving away from convex optimization

Issue #1: can't find global minimum (NP-hard)

Issue #2: minimum we find depends on opt. algo and init. Langevin Dynamics

Issue #3: generalizability

Neural Tangent Kernel

① Overview

fundamental result about fitting overparametrized deep nets

idea: for wide enough deep net GD finds a 0-error solution

concrete:

(1) depth ≥ 2 , right nonlinearities (\sim universal approx)

(2) how wide is wide enough

(3) right scaling parameters for random init

Main Ideas:

(1) if we take linear approx to a deep network, if width $>$ #examples, we can fit perfectly

(2) as we increase width, we need to move less in parameter space

→ weight doesn't change too much

① Setup

Setup

$$f(w) \triangleq \begin{bmatrix} f(x_1; w) \\ \vdots \\ f(x_n; w) \end{bmatrix} \quad y \triangleq \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

$$R(2f(w)) \triangleq \frac{1}{2} \|2f(w) - y\|^2$$

2 ~ width

$$\dot{w}(t) = -\nabla_w R(2f(w(t))) = -2j_t^T \nabla R(2f(w(t))) = -2j_t^T (2f(w(t)) - y)$$

gradient flow

$$j_t \triangleq \begin{bmatrix} \nabla f(x_1; w(t))^T \\ \vdots \\ \nabla f(x_n; w(t))^T \end{bmatrix}$$

Neim Theorem (Informal)

(1) if λ is big enough $\lambda(2f(w(t))) \in R(2f(w(t))) \subset \mathbb{C}^{n+1}$

Goal

$$(2) \|w(t) - w(0)\| \leq \frac{\sqrt{R(2f(w(0)))}}{\lambda}$$

proof strategy: introduce auxiliary flow $u(t)$ using linear approximation to f and show that $u(t)$ and $w(t)$ remain close

① look at lemma for gradient flow

$$f_0(u) \triangleq f(w(0)) + j_0(u - w(0)) \text{ linear approx (Taylor series) around initialization}$$

$$\text{gradient flow: } \dot{u}(t) = -\nabla_u R(2f_0(u(t))) = -2j_0^T \nabla R(2f_0(u(t))) = -2j_0^T (2f_0(u(t)) - y) \quad (1)$$

Jacobian doesn't depend on time now!

Lemma 1: if $j_0 j_0^T$ is full rank, then $R(2f_0(u(t)))$ goes to zero

Proof: look at change in predictions

$$\frac{d}{dt} 2f_0(u(t)) = \frac{d}{dt} \left[2(f(w(0)) + j_0(u(t) - w(0))) \right] = 2j_0 \dot{u}(t) \quad (2)$$

$$(1), (2) \Rightarrow \frac{d}{dt} 2f_0(u(t)) = -2^2 j_0 j_0^T (2f_0(u(t)) - y)$$

$$\frac{d}{dt} (\underbrace{2f_0(u(t)) - y}_{\triangleq v(t)}) = -2^2 j_0 j_0^T (2f_0(u(t)) - y).$$

$v(t) = -2^2 j_0 j_0^T v(t)$ simple ODE; solution is modute exponential

$v(t) = e^{-2^2 j_0 j_0^T t} v(0)$ converges by full rank assumption to zero

key definition: matrix $2^2 j_0 j_0^T$ is called mean tangent kernel

$K(x_i, x_j) = \nabla f(x_i; w(0))^T \cdot \nabla f(x_j; w(0))$ kernel function because can be used on new data

Note It is known that under right scaling $f(\cdot; w(t))$ converges to a Gaussian process property of architecture and random initialization; given by gradients

Note gradient flow for w is something we have seen, in disguise, for least squares

$$\text{is: } \dot{y}_{t+1} = y_t - \eta A^T \sigma_t, \quad \sigma_t \stackrel{\Delta}{=} A z_t - b$$

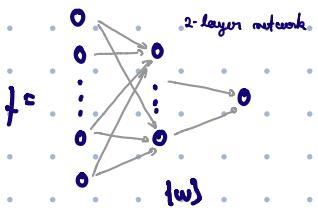
rewrite as recurrence for residual σ_t . $\sigma_{t+1} = \sigma_t - \eta A^T \sigma_t$
looks like what we had with $A \leftrightarrow J_0$

w is solving a least squares problem using the linear approx to the deep net at initialization

(IV) Intuition: why 2^d ?

2^d abstraction for why increasing 2^d (increasing width)
as $2 \rightarrow \infty$ we are becoming (in at core) so the linear approx is more accurate
as we increase width we need to move less to offset output

thought experiment $f = f_1$



\Rightarrow double the network to \hat{f}

$\{w_1, w_2\}$ & both are duplicates of w

random matrix theory to prove this for depth-2

(V) Analyze real gradient-flow $w(t)$

(1) evolution in prediction space: $\frac{d}{dt} 2f(w(t)) = 2J_e w(t) = -2^3 J_e^T \nabla R(2f(w(t)))$
2 function of time!

$$(2) motion $z(t) \stackrel{\Delta}{=} 2f(w(t)) \quad \dot{z}(t) \stackrel{\Delta}{=} 2J_e^T$$$

Lemma 2: If: $\dot{z}(t) = -Q(t) \nabla R(z(t))$

$$\lambda \stackrel{\Delta}{=} (\inf_{t \in [0, T]} \lambda_{\min}(Q(t)) > 0$$

Then: for $t \in [0, T]$, $R(z(t)) \leq R(z(0)) e^{-ct\lambda}$

Proof: 1) Gronwall's Inequality

$[a, b]$ interval $u(t) \leq p(t) u(t), \forall t \in [a, b]$

then u is upper-bounded by sol. to

$$U(t) = p(t) U(0)$$

assuming some boundary conditions
smooth
 $\Rightarrow u(t) \leq U(0) e^{-c(t-a)}$

2) look at

$$\frac{d}{dt} \frac{1}{2} \|z(t) - y\|^2 = \langle \dot{z}, z - y \rangle = \langle -Q(t) \nabla R(z(t)), z - y \rangle$$

$$= \langle -q(t) \cdot (z-y), z-y \rangle \leq -2\lambda \frac{\|z(t)-y\|^2}{2}$$

So $\frac{d}{dt} R(z(t)) \leq -2\lambda R(z(t))$. Now we can use Gronwall's Ineq.

Corollary: $\|z(t)-y\| \leq \|z(0)-y\| e^{-\lambda t}$

Next steps: show we don't go too far away from init

Notation:

$$v(t) \triangleq w(t)$$

$$g(w(t)) \triangleq f(w(t))$$

$$S(\epsilon) \triangleq 2J_\epsilon$$

Lemma 2 $\dot{v}(t) = -S(t)^T \nabla R(g(w(t)))$

$$\Theta(t) = S(t) S(t)^T, t \in [0, 2]$$

$$\lambda I \leq \Theta(t) \leq \lambda_{\max} I$$

$$\text{Thm: } \|v(t) - v(0)\| \leq \frac{\sqrt{\lambda_{\max}}}{\lambda} \|g(v(0)) - y\| \\ \leq \frac{\sqrt{2\lambda_{\max} R(g(v(0)))}}{\lambda}$$

Proof: triangle inequality + corollary

Discussion: We need $\lambda \sigma_{\min} I \leq G(\epsilon) \leq \lambda \sigma_{\max} I$ (3)

Epilogue: What random initialization puts you in nice regime?

$$x \rightarrow A_i x_i + b_i ; \text{ Let } A_i = \frac{L}{\sqrt{m_i}} w_i$$