New Insights about SGD stepsize, risk convergence, and implicit regularization

Jingfeng Wu 10/2023









Modern machine learning



AlphaGo

At their core is deep learning



Google assistant



ChatGPT

Deep learning oversimplified





Repeat for TONs of (x, y) + some tips/tricks/tuning



Current DL practice

- **Usual formula:**
- 1. one graduate student
- 2. one model
- 3. performance numbers on a few standard test sets
- 4. yay. we. rock.
- 5. one Ph.D.

https://www.slideshare.net/npinto/highperformance-computing-needs-machine-learning-and-vice-versa-nips-2011-big-learning



"Alchemy", largely trial & error

"This is Graduate Student Descent"

Current DL theory



blind men and an elephant

Would theory, from an incomplete picture, be useful?

- parameters: universal approximation, ...
- random initialization: NTK, ...
- low training error: benign overfitting, ...
- SGD: implicit bias, ...
- and many more

This talk: better theory, better trial & error

only got two grad students what to play with





the incomplete theory

its a feature not a bug can try that

Part 1. Stepsize

(Stochastic) Gradient Descent

$$w_{+} = w - \eta \cdot$$

making how much update? AKA., stepsize / learning rate?

backpropagation adaptive linear neuron



Werbos, 1974



Widow & Hoff, 1960



gradient descent

 $\nabla \ell(w)$



Cauchy, 1847

perceptron

stochastic approximation



Rosenblatt, 1958



Robbins & Monro, 1951





Optimization theory oversimplified quadratic landscape

$$\begin{split} \ell(w_{+}) &= \ell(w - \eta \cdot \nabla \ell(w)) \\ &= \ell(w) - \eta \cdot \|\nabla \ell(w)\|^{2} + \frac{\eta^{2}}{2} \cdot \nabla \ell(w)^{\top} \cdot \nabla^{2} \ell(w) \cdot \nabla \ell(w) - O(\eta^{3}) \\ &\leq \ell(w) - \eta \cdot \left(\left(1 - \frac{\eta}{2} \cdot \|\nabla^{2} \ell(w)\|_{2}\right) \cdot \|\nabla \ell(w)\|^{2} - O(\eta^{3}) \right) \end{split}$$

[descent lemma]

For small η , $\ell(w_t)$ decreases monotonically, GD works

For large η , GD does not work for quadratics





Numbers in DL (classification problem)

3-layer net + 1,000 samples from MNIST+ GD with const-stepsize



small stepsize works; large stepsize also works and works better in the long run

Limitations of the old picture

DL classification (cross-entropy)



old theory picture predicts green curve, but fails predicting red curve



opt theory (quadratic)



Logistic regression, revisited

[WBL'23] In the NN training (minimizing a cross-entropy loss), if

- the model is linear, i.e., $f(x) = x^{\top}w$,
- the datasets can be perfectly fitted,

Then GD, with **any constant stepsize**, minimizes the fitting error.

[open/ongoing]

Benefits of large stepsize?

Wu, Jingfeng, Vladimir Braverman, and Jason D. Lee. "Implie Stability." NeurIPS 2023



Wu, Jingfeng, Vladimir Braverman, and Jason D. Lee. "Implicit Bias of Gradient Descent for Logistic Regression at the Edge of

Rethinking optimization theory

- in DL (classification), Taylor approx. is limited (b/c its local)
- "valley" is more stable than "quadratic basin"
- descent lemma isn't a golden rule (more likely a bad rule)





Part 2. Implicit Regularization

Statistical learning theory oversimplified

 $R(w) - R(w^*) = R(w) - \hat{R}_n(w)$

 \leq Empirical +

1 is (always) a good advice, 2 is debatable

 $\sup_{w \in \mathscr{U}} |R(w) - \hat{R}_n(w)| \lesssim \sqrt{\frac{\log |\mathscr{H}|}{n}}$ $w \in \mathcal{H}$

2. **#Param** (uniformly) controls complexity

$$+\hat{R}_{n}(w) - \hat{R}_{n}(w^{*}) + \hat{R}_{n}(w^{*}) - R(w^{*})$$

$$2 \cdot \sup_{w \in \mathcal{H}} |R(w) - \hat{R}_{n}(w)|^{4}$$

1. Generalization: empirical fit vs. complexity control



What controls complexity in DL?

- 1. hypothesis class size
 - #params
- 2. explicit regularization
 - weight decay, dropout,...
- 3. implicit regularization
 - SGD, a simple learning rule

How does SGD regularize the model?

Wu, Jingfeng, Difan Zou, Vladimir Braverman, and Quanquan Gu. "Direction Matters: On the Implicit Bias of Stochastic Gradient Descent with Moderate Learning Rate." ICLR 2021

far from overfitting



NN with **11,330** parameters

fitting **2,000** samples from FashionMNIST

without explicit regularization



Linear / ReLU regression, revisited minimize $R(\mathbf{w}) = \mathbb{E}(\phi(\mathbf{x}^{\mathsf{T}}\mathbf{w}) - y)^2$, $\mathbf{w} \in \mathbb{R}^d$ $\phi = id \text{ or } max\{\cdot, 0\}$ with n iid samples: $(\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_n, y_n)$ **Distributional assumption** (simplified)

 $y = \phi(\mathbf{x}^{\mathsf{T}} \mathbf{w}_*) + \mathcal{N}(0,1), \quad \|\mathbf{w}_*\|_2 \le 1, \quad \mathbf{x} \sim \mathcal{N}(0,\mathbf{H})$

SGD / Perceptron

 $\mathbf{w}_t = \mathbf{w}_{t-1} - \eta_t \cdot \left(\phi(\mathbf{x}_t^{\mathsf{T}} \mathbf{w}_{t-1}) - y_t \right) \cdot \mathbf{x}_t, \quad t = 1, \dots, n \qquad \underset{0.1\eta}{0.1\eta}$



Implicit complexity control

[ZWBGK'21, WZBGK'22, WZCBGK'23] $Bias + \frac{DIM}{m} \lesssim \mathbb{E}R(\mathbf{w}_n) - \min R \lesssim Bias + \frac{DIM}{m}$ n

$DIM \leq d$, a function of n, η , and H, controls the complexity

- regression." COLT 2021.
- stepsize for overparameterized linear regression." ICML 2022
- High-Dimensional Single ReLU Neuron." ICML 2023



• Zou, Difan, Jingfeng Wu, Vladimir Braverman, Quanquan Gu, and Sham Kakade. "Benign overfitting of constant-stepsize sgd for linear

• Wu, Jingfeng, Difan Zou, Vladimir Braverman, Quanquan Gu, and Sham Kakade. "Last iterate risk bounds of sgd with decaying

• Wu, Jingfeng, Difan Zou, Zixiang Chen, Vladimir Braverman, Quanquan Gu, and Sham M. Kakade. "Finite-Sample Analysis of Learning



Implicit complexity control $\mathsf{DIM} := \# \left\{ \lambda_i \ge \frac{1}{n\eta} \right\} + n^2 \eta^2 \cdot \sum_{\lambda_i < \frac{1}{m}} \lambda_i^2 \quad (\lambda_i)_{i \ge 1} \text{ denote eigenvalues of } \mathbf{H} = \mathbb{E}[\mathbf{x}\mathbf{x}^\top]$



SGD exploits a "low-dim" structure in data

DIM is small when eigenvalues decay fast



Implicit vs. explicit regularization

[ZWBGFK'21] Fix a set of "natural" least square problems (opt is not hard & label is noisy).

Assume an oracle tunes both SGD and ridge to their best on each problem instance. Let m and n be their sample complexity to attain the same rate.

- For every problem instance:
- There **exist** a problem instance:

tuned SGD inflates **at most poly-log** more samples than tuned ridge tuned ridge can inflate **poly** more samples than tuned SGD

tune SGD before tuning weight decay :)

Zou, Difan, Jingfeng Wu, Vladimir Braverman, Quanquan Gu, Dean P. Foster, and Sham Kakade. "The benefits of implicit regularization from sgd in least squares problems." NeurIPS 2021

$$m \lesssim n \log^2(n)$$

$$n \ge m^2/\log^4(m)$$

Rethinking statistical learning theory

- in DL, complexity isn't controlled by #params
- complexity is largely controlled by simple learning rules
- stepsize balances empirical fit and complexity control



problem instance



Conclusions





- optimization and statistical learning theory are still insightful
- but need revisions in deep learning
- large stepsize, implicit regularization
- more effective trial & error if things are understood better

