The Implicit Regularization of SGD in Least Square Problems and Beyond

Jingfeng Wu, 01/2023

DL 101: How to train DNNs on CIFAR-10?

Bag of tricks

- Data augmentation
- Weight decaying
- Dropout

- - -

Batch normalization

SGD (LR decaying + early stopping)

airplane automobile bird cat deer dog frog horse ship

truck





Algorithms induce regularization

• "Unregularized":

- "Explicit" regularization: weight decaying / ridge $w \leftarrow \arg\min L(w) + \lambda \|w\|_2^2$
- "Implicit" regularization: SGD
 - $w \leftarrow SGD(\eta; dataset)$

 $w \leftarrow \arg\min L(w)$

SGD in **Practice**

- Overparameterized model => Tons of ERM & some may not generalize SGD: small batch + large initial LR + LR decaying (early stopping)
- SGD solution generalizes well

Algo implicitly imposes regularization! But how?

Wu, Jingfeng, Wenqing Hu, Haoyi Xiong, Jun Huan, Vladimir Braverman, and Zhanxing Zhu. "On the noisy gradient descent that generalizes as sgd." In International Conference on Machine Learning, pp. 10367-10376. PMLR, 2020.

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \cdot \nabla \mathscr{C}(\mathbf{x}_i, y_i; \mathbf{w})$$



Q: For a new problem, which trick to try/tune first?

SGD or WD?

Understand SGD



Problem simplification









- # param > 10^{7}
- ReLU non-linearity
- Layer structure



Maybe...

Start with modeling high-dim?

Head up: + ReLU shortly



A high-dim linear regression

• *d* dimensional regression

• A small training set (iid), n < d

 $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$

• Test error

 $\Delta(\mathbf{w}) := \|\mathbf{w} - \mathbf{w}^*\|_{\mathbf{H}}^2$

$y = \mathbf{x}^{\top} \mathbf{w}_{*} + \mathcal{N}(0,1), \|\mathbf{w}_{*}\|_{2} \le 1, \mathbf{x} \sim \mathcal{N}(0,\mathbf{H})$

Wait, is this even possible to solve? In general, no

[Classical Result]

[Intuition]

- probe each dim induces a unit uncertainty
- d important directions to probe \bullet

When $\mathbf{H} = \mathbf{I}_d$, any reasonable algorithm suffers $\Delta \geq \frac{d}{-} \geq \Omega(1)$

But why we can solve them in practice?



In practice, H only has a few large eigenvalues. $\mathbf{H} \ll \mathbf{I}$

Yang, Rubing, Jialin Mao, and Pratik Chaudhari. "Does the data induce capacity control in deep learning?." In International Conference on Machine Learning, pp. 25166-25197. PMLR, 2022.



Q: How SGD performs when $H \ll I$?

Algorithm simplifaction

One-pass SGD (i.e., early stopping)

$\mathbf{w}_t = \mathbf{w}_{t-1} - \eta_t \cdot (\mathbf{x}_t^{\top} \mathbf{w}_{t-1} - y_t) \cdot \mathbf{x}_t, \quad t = 1, \dots, n$







problem instance

$$\begin{array}{l} \frac{1}{2} \cdot \text{DIM} \leq \Delta \leq \frac{\text{const2}}{n} \cdot \text{DIM} \\ \text{DIM} := \#\left\{ \left(\lambda_i\right) \geq \frac{1}{\eta n} \right\} + \eta^2 n^2 \sum_{\lambda_i < \frac{1}{\eta n}} \\ \text{eigenvalues of } \mathbf{H} = \mathbb{E}[\mathbf{x}\mathbf{x}^\top] \end{array}$$

- Zou, Difan, Jingfeng Wu, Vladimir Braverman, Quanquan Gu, and Sham Kakade. "Benign overfitting of constant-stepsize sgd for linear regression." In Conference on Learning Theory, pp. 4633-4635. PMLR, 2021.
- Wu, Jingfeng, Difan Zou, Vladimir Braverman, Quanquan Gu, and Sham M. Kakade. "Last Iterate Risk Bounds of SGD with Decaying Stepsize for Overparameterized Linear Regression." International Conference on Machine Learning, 2022.





SGD can solve well-structured high-dim problems

SGD or WD?

SGD vs. WD

• Ridge regularization (WD)

$$\mathbf{w} = \arg\min \sum_{t=1}^{n} |$$

One-pass SGD (i.e., early stopping)

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \eta_t \cdot (\mathbf{x}_t^{\mathsf{T}} \mathbf{w}_{t-1} - y_t) \cdot \mathbf{x}_t, \quad t = 1, \dots, N$$



$\|\mathbf{x}_n^{\mathsf{T}}\mathbf{w} - y_n\|_2^2 + \lambda \cdot \|\mathbf{w}\|_2^2$

SGD vs. WD



problem instance

- Tsigler, Alexander, and Peter L. Bartlett. "Benign overfitting in ridge regression." arXiv preprint arXiv:2009.14286 (2020).
- Information Processing Systems 34 (2021): 5456-5468.

Ridge could be bad 😕 SGD is alway nearly good \bigcirc

Tune both SGD and ridge to their best.

Let *m* and *n* be their sample complexity

For each problem in set: lacksquare

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

There is a problem in the set: •

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

• 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." Advances in Neural



Tuning SGD is probably more worthy than tuning WD



Beyond Linearity



A high-dim ReLU regression

$y = \operatorname{ReLU}(\mathbf{x}^{\mathsf{T}}\mathbf{w}_{*}) + \mathcal{N}(0,1)$

Non-convex -> very hard

• SGD no longer a good idea

• But, can fix it!



A loss landscape example

Non-convex

• SGD, initialized from left, will stuck

 Exact-gradient based updates may not be good



ML101: Perceptron $L(\mathbf{w}) = \sum_{(\mathbf{x}, y)} (\mathbf{F})$

• SGD

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \eta_t \cdot (\text{ReLU}(\mathbf{x}_t^{\top} \mathbf{w}_{t-1}) - y_t)$$

Perceptron (iteratively error correction)

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \eta_t \cdot (\text{ReLU}(\mathbf{x}_t^\top \mathbf{w}_{t-1}))$$

 $L(\mathbf{w}) = \sum \left(\text{ReLU}(\mathbf{x}^{\top}\mathbf{w}_{*}) - y \right)^{2}$





problem instance

$$\text{DIM} := \#\left\{ \lambda_i \ge \frac{1}{\eta n} \right\} + \eta^2 n^2 \sum_{\lambda_i < \frac{1}{\eta n}}$$

How to extend perceptron to DL?





Better theory

• High-dim

. . .

- Non-convex
- Average cases



Experiments

Better practice

- Algo. choice
 - SGD > WD
- Algo. design
 - Explore structure
 - Adjust gradient?





Take Home & Acknowledgement

- Algorithmic regularization
- **SGD** > **WD**
- ReLU: perceptron > SGD



Vova Bravermen

Dean P. Foster

Quanquan Gu

Sham M. Kakade

Difan Zou