Machine Learning: Generalization 4

Hao Tang

March 24, 2025

L₂ Regularization

- The term $\frac{\lambda}{2} \|w\|_2^2$ is called an L_2 regularizer.
- It is also known as weight decay.
- The expression

$$L_{\mathcal{S}}(w) + \frac{\lambda}{2} \|w\|_2^2 \tag{1}$$

is the Lagrangian of

$$\min_{w} L_{S}(w)$$
(2)
s.t. $||w||_{2} \leq B$ (3)

Generalization bound for bounded linear classifier

• With probability $1 - \delta$, for all $h \in \mathcal{H}$,

$$L_{\mathcal{D}}(h) \leq L_{\mathcal{S}}(h) + \sqrt{\frac{r^2 B^2}{n}} + 3\sqrt{\frac{\log(2/\delta)}{2n}},$$

where $\|x\|_2 \leq r$ for any $x \in S$ and $\mathcal{H} = \{x \mapsto w^\top x : \|w\|_2 \leq B\}.$

3/31

(4)

- A learning algorithm is **stable** if the learned program does not change much in performance when we change the data set slightly.
- The slight change in data set is by swapping out a data point.

$$S = \{(x_1, y_1), \dots, (x_i, y_i), \dots, (x_n, y_n)\}$$
(5)

$$S^{i} = \{(x_{1}, y_{1}), \dots, (x', y'), \dots, (x_{n}, y_{n})\}$$
 (6)

• A learning algorithm is stable is A(S) and $A(S^{i})$ is "similar," or

$$\ell(A(S)(x), y) - \ell(A(S^{i})(x), y)$$
(7)

is small.¹

¹Recall that $L_S(h) = \frac{1}{n} \sum_{i=1}^n \ell(h(x_i), y_i).$

• Stable learning algorithms don't overfit.

 $\mathbb{E}_{S \sim \mathcal{D}^n}[L_{\mathcal{D}}(\mathcal{A}(S)) - L_S(\mathcal{A}(S))] = \mathbb{E}_{\substack{i \sim U(n) \\ S \sim \mathcal{D}^n \\ (x,y) \sim \mathcal{D}}} [\ell(\mathcal{A}(S^i)(x_i), y_i) - \ell(\mathcal{A}(S)(x_i), y_i)]$ (8)

• Proof

 $\mathbb{E}_{\mathcal{S}}[\mathcal{L}_{\mathcal{D}}(\mathcal{A}(\mathcal{S}))] = \mathbb{E}_{\mathcal{S}}[\mathbb{E}_{(x,y)\sim\mathcal{D}}[\ell(\mathcal{A}(\mathcal{S})(x), y)]] = \mathbb{E}_{\mathcal{S}}[\mathbb{E}_{(x,y)\sim\mathcal{D}}[\ell(\mathcal{A}(\mathcal{S}^{i})(x_{i}), y_{i})]]$ (9)

$$\mathbb{E}_{S}[L_{S}(A(S))] = \mathbb{E}_{S}[\mathbb{E}_{i \sim U(n)}[\ell(A(S)(x_{i}), y_{i})]]$$
(10)

• If ℓ is convex and ρ -Lipschitz² and $A(S) = \operatorname{argmin}_{h \in \{w: x \mapsto w^{\top}x\}} [L_S(w) + \lambda ||w||_2^2]$, then

$$\|A(S^{i}) - A(S)\|_{2} \leq \frac{2\rho}{\lambda n}.$$
(11)

• In the end, we have

$$\mathbb{E}_{\mathcal{S}\sim\mathcal{D}^n}\left[L_\mathcal{D}(\mathcal{A}(\mathcal{S}))-L_\mathcal{S}(\mathcal{A}(\mathcal{S}))
ight]\leq rac{2
ho^2}{\lambda n}.$$

²A function f is ρ -Lipschitz if $|f(x) - f(y)| \le \rho ||x - y||_2$ for any x and y.

(12)

- Minimizing L₂ regularized convex and Lipschitz functions is a stable learning algorithm.
- Stable learning algorithms don't overfit.

Computational

Statistical

Computational

Statistical

Runtime

Samples

Computational

Statistical

Runtime

Samples

How many steps do we need?

How many samples do we need?

Computational

Statistical

Runtime

Samples

How many steps do we need?

Polynomial number of steps

How many samples do we need?

Polynomial number of samples

VC dimension of a sine function



- For every ε > 0, given any Lipschitz function f : [-1, 1]^d → [-1, 1], there is a network g such that |g(x) f(x)| ≤ ε for any x.
- The number of nodes needed to achieve this is $O(2^d)$.







Can we approximate a sine function?



- Polynomials are universal approximators.
- Decision trees are universal approximators.
- Gaussian mixture models are universal approximators.
- Universal approximation does not explain why neural networks are so "special."

Depth separation

- There exists functions which can be approximated with small depth 3 networks, but cannot be approximated with depth 2 networks without using $O(2^d)$ nodes.
- Functions to show these results tend to oscillate a lot.
- Some believe the results are pathological and do not happen in practice.

- What can be implemented with polynomial number of of nodes?
- Any Turing machine that runs in T operations can be implemented with a neural network of depth O(T) with a total $O(T^2)$ nodes.
- Recall that VC dimension of neural networks is $O(|E| \log |E|)$, where E is the number of edges in the network.

• Training a 2-layer 3-node neural network is NP-complete.

- Training a 2-layer 3-node neural network is NP-complete.
- The proof converts instances of an NP-complete problem into data points.
- If we can minimize the loss of the training set, we solve the NP-complete problem.

- Training a 2-layer 3-node neural network is NP-complete.
- The proof converts instances of an NP-complete problem into data points.
- If we can minimize the loss of the training set, we solve the NP-complete problem.
- Maybe we don't need to solve this exactly?

- Approximating ERM is NP hard.
- The loss is not necessarily convex.
- ERM is hard for neural networks.

Optimizing neural networks on random labels





- Overparameterization means using a lot more nodes than the number of points.
- Overparameterization helps optimization.

- Overparameterization means using a lot more nodes than the number of points.
- Overparameterization helps optimization.
- Wouldn't the model just memorize the training set?
- Wouldn't the hypothesis class be too large to have good generalization error?





Image credit: (Neyshabur et al., 2014)

- Fitting a data set to training error zero is called interpolation.
- Why doesn't interpolation overfit?









Overfitting



Image credit: (Mallinar et al., 2022)

Sharp and flat minima



Image credit: (Foret et al., 2021)

The grand goal

• Coming back to regularized ERM, it considers both the training error and the capacity during optimization.

$$\min_{w} \quad L_{S}(w) + \frac{\lambda}{2} \|w\|_{2}^{2}$$
(13)

• If we know something that controls the capacity, we should optimize it.

In practice

- Always start with the training error.
- Always start with ERM.
- Why is the training error not close to zero?
- Regularize