Techniques for Improving Generalization





L2 regularization: $\frac{\lambda}{2} \|\theta\|_2^2$

Implementation: $\theta \leftarrow (1 - \eta \lambda)\theta - \eta \nabla f(\theta)$

Dropout

Intuition: randomly cut off some connections and neurons.

Training: for each input, at each iteration, randomly "turn off" each neuron with a probability $1 - \alpha$

- Change a neuron to 0 by sampling a Bernoulli variable.
- Gradient only propogatd from non-zero neurons.





Dropout

Dropout changes the scale of the output neuron:

- $y = \text{Dropout}(\sigma(WX))$
- $\mathbb{E}[y] = \alpha \mathbb{E}[\sigma(Wx)]$

Test time: $y = \alpha \sigma(Wx)$ to match the scale





Understanding Dropout

- Dropout forces the neural network to learn redundant patterns.
- Dropout can be viewed as an implicit L2 regularizer (Wager, Wang, Liang '13).





Early Stopping

- Continue training may lead to overfitting.
- Track performance on a held-out validation set.
- Theory: for linear models, equivalent to L2 regularization.



Data Augmentation

Depend on data types.

Computer vision: rotation, stretching, flipping, etc



CocaColaZero1_1.png



CocaColaZero1_5.png



CocaColaZero1_2.png



CocaColaZero1_6.png



CocaColaZero1_3.png



CocaColaZero1_7.png



CocaColaZero1_4.png



CocaColaZero1_8.png

Mixup data augmentation

- $\hat{x} = \lambda x_i + (1 \lambda) x_j$
- $\hat{y} = \lambda y_i + (1 \lambda) y_j$
- $\lambda \sim \text{Beta}(0.2)$



Data Augmentation

Depend on data types.

- Natural language processing:
- Synonym replacement
 - This article will focus on summarizing data augmentation in NLP.
 - This write-up will focus on summarizing data augmentation in NLP.
- Back translation: translate the text data to some language and then translate back
 - I have no time. -> 我没有时间. -> I do not have time.

Learning rate scheduling

Start with large learning rate. After some epochs, use small learning rate.



Learning rate scheduling

Start with large learning rate. After some epochs, use small learning rate.

Theory:

- Linear model / Kernel: large learning rate first learns eigenvectors with large eigenvalues (Nakkiran, '20).
- Representation learning (Li et al., '19)



Normalizations

- Batch normalization (loffe & Szegedy, '15)
- Layer normalization (Ba, Kiros, Hinton, '16)
- Weight normalization (Salimans, Kingma, '16)
- Instant normalization (Ulyanov, Vedaldi, Lempitsky, '16)
- Group normalization (Wu & He, '18)

Generalization Theory for Deep Learning



Basic version: finite hypothesis class

Finite hypothesis class: with probability $1 - \delta$ over the choice of a training set of size *n*, for a bounded loss ℓ , we have

$$\sup_{f \in \mathscr{F}} \left| \frac{1}{n} \sum_{i=1}^{n} \ell(f(x_i), y_i) - \mathbb{E}_{(x, y) \sim D} \left[\ell(f(x), y) \right] \right| = O\left(\sqrt{\frac{\log|\mathscr{F}| + \log 1/\delta}{n}}\right)$$

VC-Dimension

Motivation: Do we need to consider **every** classifier in \mathcal{F} ? Intuitively, **pattern of classifications** on the training set should suffice. (Two predictors that predict identically on the training set should generalize similarly).

Let
$$\mathscr{F} = \{f : \mathbb{R}^d \to \{+1, -1\}\}$$
 be a class of binary classifiers.

The growth function $\Pi_{\mathscr{F}} : \mathbb{N} \to \mathbb{F}$ is defined as:

$$\Pi_{\mathscr{F}}(m) = \max_{(x_1, x_2, \dots, x_m)} \left| \left\{ (f(x_1), f(x_2), \dots, f(x_m)) \mid f \in \mathscr{F} \right\} \right|.$$

The VC dimension of \mathscr{F} is defined as: VCdim $(\mathscr{F}) = \max\{m : \Pi_{\mathscr{F}}(m) = 2^m\}$.

VC-dimension Generalization bound

Theorem (Vapnik-Chervonenkis): with probability $1 - \delta$ over the choice of a training set, for a bounded loss ℓ , we have

$$\sup_{f \in \mathscr{F}} \left| \frac{1}{n} \sum_{i=1}^{n} \mathscr{C}(f(x_i), y_i) - \mathbb{E}_{(x, y) \sim D} \left[\mathscr{C}(f(x), y) \right] \right| = O\left(\sqrt{\frac{\mathsf{VCdim}(\mathscr{F})\log n + \log 1/\delta}{n}}\right)$$

Examples:

- Linear functions: VC-dim = O(dimension)
- Neural network: VC-dimension of fully-connected net with width W and H layers is $\Theta(WH)$ (Bartlett et al., '17).

Problems with VC-dimension bound

- 1. In over-parameterized regime, bound >> 1.
- 2. Cannot explain the random noise phenomenon:
 - Neural networks that fit random labels and that fit true labels have the same VC-dimension.



PAC Bayesian Generalization Bounds

Setup: Let *P* be a prior over function in class \mathcal{F} , let *Q* be the posterior (after algorithm's training).

Theorem: with probability $1 - \delta$ over the choice of a training set, for a bounded loss ℓ , we have

$$\sup_{f \in \mathscr{F}} \left| \frac{1}{n} \sum_{i=1}^{n} \mathscr{\ell}(f(x_i), y_i) - \mathbb{E}_{(x, y) \sim D} \left[\mathscr{\ell}(f(x), y) \right] \right| = O\left(\sqrt{\frac{KL(Q \mid |P) + \log 1/\delta}{n}}\right)$$

Rademacher Complexity

Intuition: how well can a classifier class fit random noise?

(Empirical) Rademacher complexity: For a training set $S = \{x_1, x_2, ..., x_n\}$, and a class \mathscr{F} , denote: $\hat{R}_n(S) = \mathbb{E}_{\sigma} \sup_{f \in \mathscr{F}} \sum_{i=1}^n \sigma_i f(x_i)$. where $\sigma_i \sim \text{Unif}\{+1, -1\}$ (Rademacher R.V.).

(Population) Rademacher complexity:

$$R_n = \mathbb{E}_S \left[\hat{R}_n(s) \right].$$

Rademacher Complexity Generalization Bound

Theorem: with probability $1 - \delta$ over the choice of a training set, for a bounded loss ℓ , we have

$$\sup_{f \in \mathscr{F}} \left| \frac{1}{n} \sum_{i=1}^{n} \mathscr{E}(f(x_i), y_i) - \mathbb{E}_{(x, y) \sim D} \left[\mathscr{E}(f(x), y) \right] \right| = O\left(\frac{\hat{R}_n}{n} + \frac{\log 1/\delta}{n}\right)$$

and

$$\sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^{n} \ell(f(x_i), y_i) - \mathbb{E}_{(x, y) \sim D} \left[\ell(f(x), y) \right] \right| = O\left(\frac{R_n}{n} + \frac{\log 1/\delta}{n}\right)$$

Kernel generalization bound

Use Rademacher complexity theory, we can obtain a generalization bound $O(\sqrt{y^{\top}(H^*)^{-1}y/n})$ where $y \in \mathbb{R}^n$ are n labels, and $H^* \in \mathbb{R}^{n \times n}$ is the kernel (e.g., NTK) matrix.



Norm-based Rademacher complexity bound

Theorem: If the activation function is σ is ρ -Lipschitz. Let $\mathscr{F} = \{x \mapsto W_{H+1}\sigma(W_h\sigma(\cdots\sigma(W_1x)\cdots), \|W_h^T\|_{1,\infty} \leq B \forall h \in [H]\}$ then $R_n(\mathscr{S}) \leq \|X^\top\|_{2,\infty}(2\rho B)^{H+1}\sqrt{2\ln d}$ where $X = [x_1, \dots, x_n] \in \mathbb{R}^{d \times n}$ is the input data matrix.