

HW 2 update

proper presentation

Techniques for Improving Generalization

W

Weight Decay

$$\min_{\theta} f(\theta) + \underbrace{\frac{\lambda}{2} \|\theta\|_2^2}$$

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

decay

Implementation: $\theta \leftarrow \underbrace{(1 - \eta\lambda)}_{\text{decay}} \theta - \eta \nabla f(\theta)$

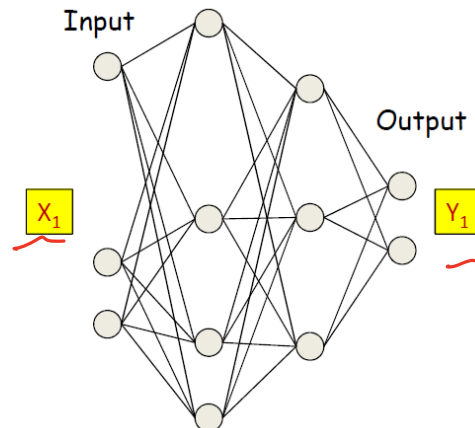
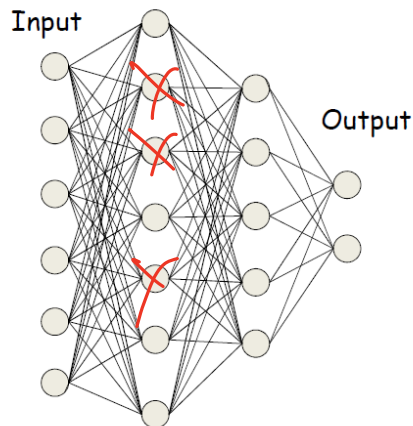
η : step size

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 propagated from non-zero neurons.

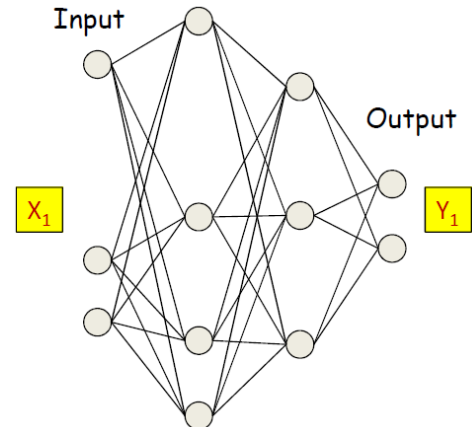
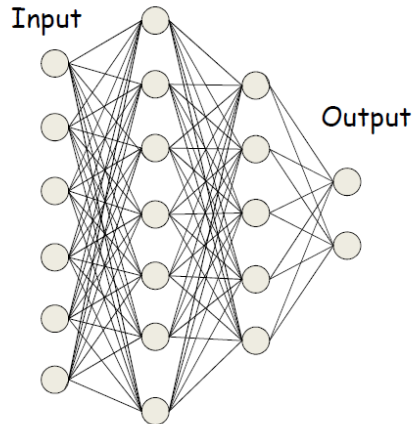


Dropout

Dropout changes the scale of the output neuron:

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

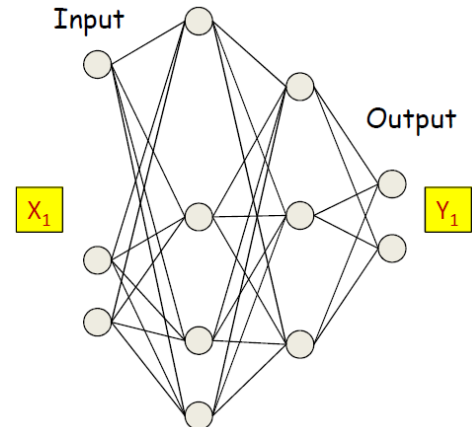
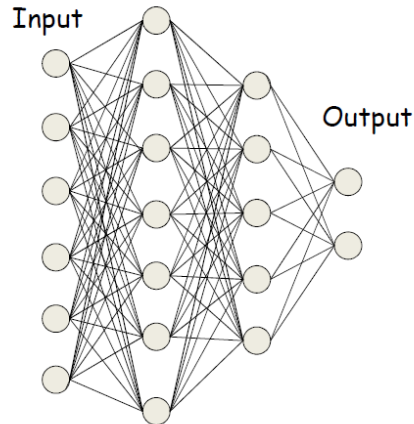
Test time: $y = \alpha \sigma(Wx)$ to match the scale



Understanding Dropout

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

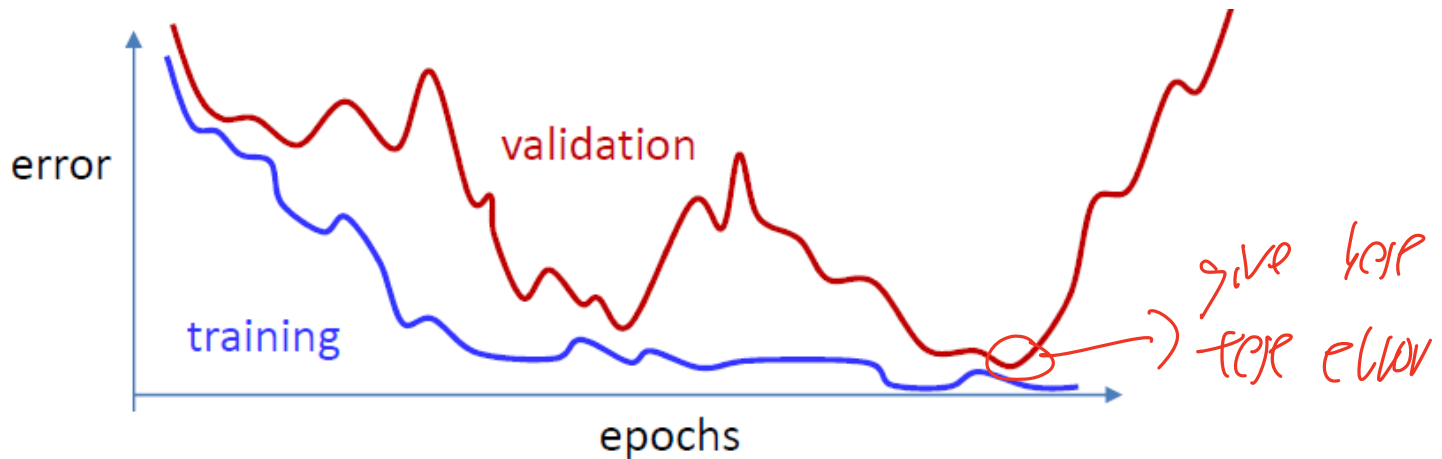
Linear, in expertly



Early Stopping

view training time as a hyper-parameter

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



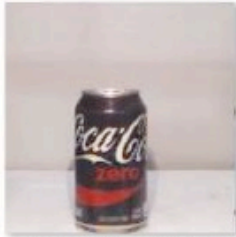
Data Augmentation

$$\sqrt{\frac{\text{comp}(f)}{n}}$$

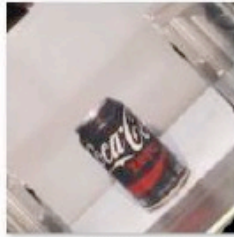
Depend on data types.

Computer vision: rotation, stretching, flipping, etc

not i.i.d.



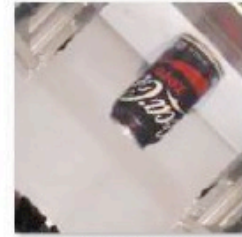
CocaColaZero1_1.png



CocaColaZero1_2.png



CocaColaZero1_3.png



CocaColaZero1_4.png



CocaColaZero1_5.png



CocaColaZero1_6.png



CocaColaZero1_7.png



CocaColaZero1_8.png



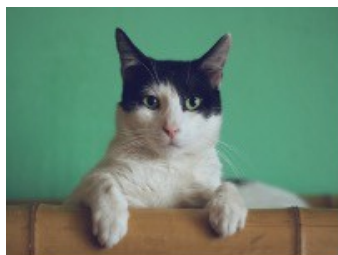
Mixup data augmentation

- $\hat{x} = \lambda x_i + (1 - \lambda)x_j$
- $\hat{y} = \lambda y_i + (1 - \lambda)y_j$
- $\lambda \sim \mathbf{Beta}(0.2) \leftarrow [0, 1]$

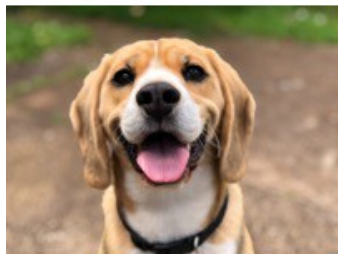
(x_i, y_i)

(x_j, y_j)

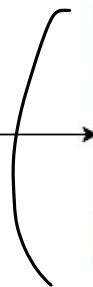
add (\hat{x}, \hat{y}) to dataset



$\begin{pmatrix} 1 \\ 0 \end{pmatrix}$



$\begin{pmatrix} 0 \\ 1 \end{pmatrix}$



$\begin{pmatrix} \lambda \\ 1-\lambda \end{pmatrix}$

Data Augmentation

Depend on data types.

sentiment analysis
sentence $\rightarrow (+, -)$

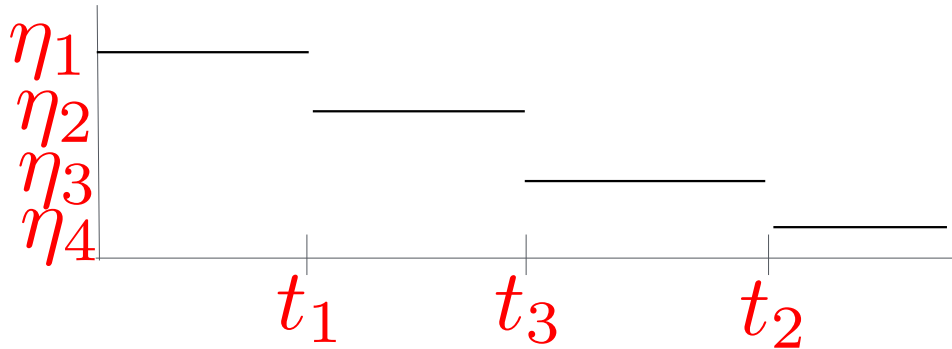
Natural language processing:

- 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.* \rightarrow 我没有时间. \rightarrow *I do not have time.*

Learning rate scheduling

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

Learning rate schedule

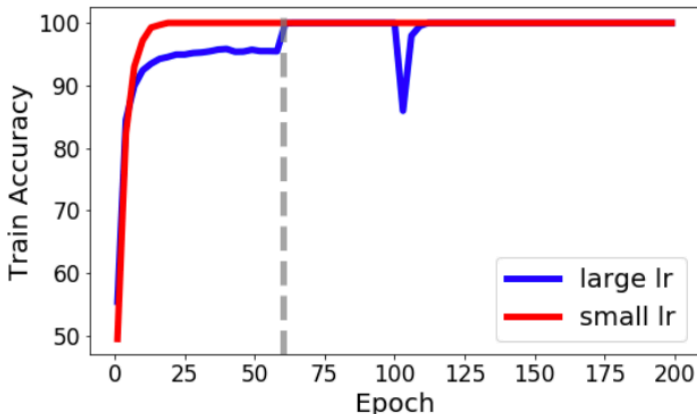


Learning rate scheduling

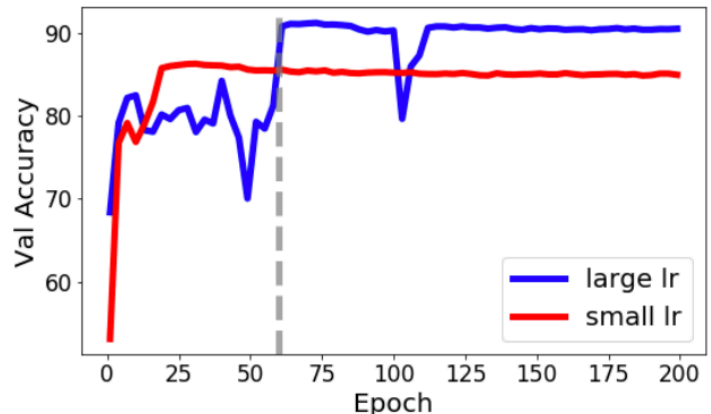
Start with large learning rate. After some epochs, use small learning rate.

Theory:

- Linear model / Kernel: large learning rate first learns eigenvectors with large eigenvalues (Nakkiran, '20).
- Representation learning (Li et al., '19)
learn coarse feature \rightarrow better generalization



Train



Validation

Normalizations

- Batch normalization (Ioffe & 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

$$\sqrt{\frac{\text{comp}(\mathcal{F})}{n}}$$

W

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

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

• For a fixed $f \in \mathcal{F}$, by Hoeffding inequality, w.p. $1 - \delta$

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

• union bound: event-1, ..., event-m

$$P \left(\bigcup_{j=1}^m \text{event-}j \right) \leq \sum_{j=1}^m P(\text{event-}j)$$

Let event- j be the event the j th function in \mathcal{F} s.t.
gen-error- $j > \sqrt{\frac{\log |\mathcal{F}| + \log 1/\delta}{n}}$, $P(\text{event-}j) \leq \delta$

$$P(\exists f \in \mathcal{F}, \text{gen-error}_f > \sqrt{\frac{\log(|\mathcal{F}|)}{n}})$$

$$= P\left(\bigcup_j \text{event}_j\right)$$

$$\leq |\mathcal{F}| \cdot \delta$$

$$\text{let } \delta' = |\mathcal{F}| \cdot \delta$$

$$\Rightarrow \sqrt{\frac{\log(|\mathcal{F}|) + \log(1/\delta)}{n}}$$

$$z_1, \dots, z_M$$

$$\left| \hat{z}_n - E(\bar{z}) \right| \approx O\left(\frac{1}{\sqrt{n}}\right)$$

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

The **growth function** $\Pi_{\mathcal{F}} : \mathbb{N} \rightarrow \mathbb{F}$ is defined as: 2^m

$$\Pi_{\mathcal{F}}(m) = \max_{(x_1, x_2, \dots, x_m)} \left| \left\{ \underbrace{(f(x_1), f(x_2), \dots, f(x_m))}_{\text{pattern}} \mid f \in \mathcal{F} \right\} \right|.$$

The **VC dimension** of \mathcal{F} is defined as:

$$\text{VCdim}(\mathcal{F}) = \max\{m : \Pi_{\mathcal{F}}(m) = 2^m\}.$$

2^d lower

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

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

Examples:

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

Problems with VC-dimension bound

$W(H) \gg n$

$\frac{VCdim}{n}$

need
data-depend

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

training:

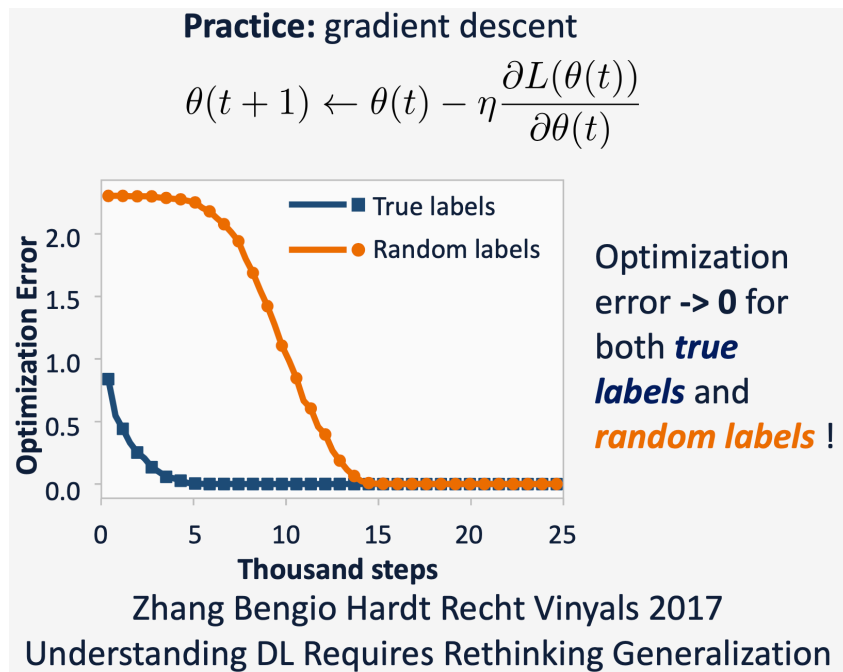
both true,
random $\rightarrow 0$

testing:

true: small
random: $\frac{1}{2}$

gen error:

true: small
random: $\frac{1}{2}$



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

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

(cannot choose (\cdot) after training)

Rademacher Complexity ± 1

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

(Empirical) **Rademacher complexity**: For a training set $S = \{x_1, x_2, \dots, x_n\}$, and a class \mathcal{F} , denote:

$$\hat{R}_n(S) = \mathbb{E}_\sigma \sup_{f \in \mathcal{F}} \sum_{i=1}^n \sigma_i f(x_i) .$$

$\langle \sigma_i, f(x_i) \rangle$
 $f(x_i)$

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

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

$\hat{R}_n = O(\sqrt{n})$

and

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

Kernel generalization bound

$y \in \mathbb{R}^n$ in training

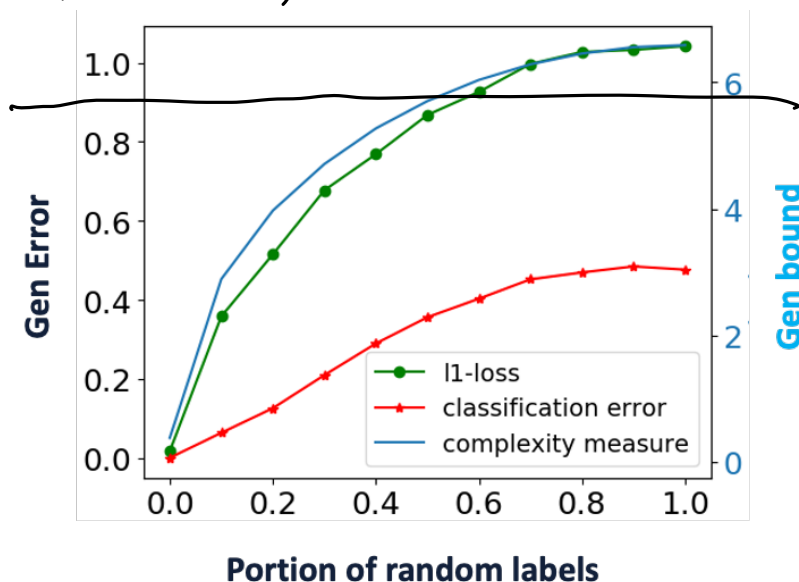
H^* depend on

(x_1, \dots, x_n)

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.

$$[f]_{ij}^\top = K(x_i, x_j)$$

✓



kernel $\rightarrow \phi$
predictor $\phi^\top w$

min $\|w\|_2$

s.t. $y_i = \phi(x_i)^\top w$

Prod: $y^\top (H^*)^{-1} y$

Norm-based Rademacher complexity bound

ReLU: $\rho=1$

Theorem: If the activation function is σ is ρ -Lipschitz. Let

$$\mathcal{F} = \{x \mapsto W_{H+1}\sigma(W_H\sigma(\cdots\sigma(W_1x)\cdots), \|W_h^T\|_{1,\infty} \leq B \forall h \in [H]\}$$

then $R_n(\mathcal{F}) \leq \|X^T\|_{2,\infty} (2\rho B)^{H+1} \sqrt{2 \ln d}$ where

$X = [x_1, \dots, x_n] \in \mathbb{R}^{d \times n}$ is the input data matrix.

$$x \in \mathbb{R}^d, w_1 \in \mathbb{R}^{m \times d}, w_2, \dots, w_H \in \mathbb{R}^{m \times m}, w_{H+1} \in \mathbb{R}^{m \times m}$$

$$\|w_1^T\|_{1,\infty} = \max_{i=1,\dots,m} \|w_1^T(\cdot, i)\|_1 = \max_{i=1,\dots,m} \sum_{k=1}^d |w_1^T(k, i)|$$

$$\|w_h^T\|_{1,\infty} = \max_{i=1,\dots,m} \|w_h^T(\cdot, i)\|_1$$

$$\|X^T\|_{2,\infty} = \max_{i=1,\dots,d} \|X^T(\cdot, i)\|_2$$