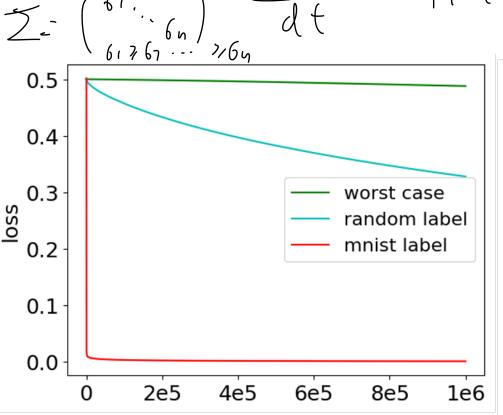
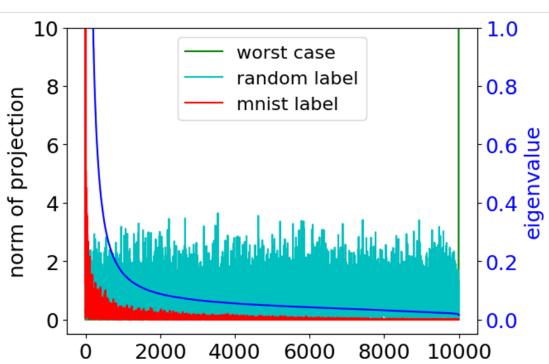
# **Neural Tangent Kernel**



What determines the convergence rate?  $\int_{-\infty}^{\infty} \frac{du(t)}{dt} = -H^*(u(t)-y)$   $\int_{-\infty}^{\infty} \frac{du(t)}{dt} = -H^*(u(t)-y)$ 





**Convergence Rate** 

**Projections** 

worst (ase;
$$y = (\cdot Un)$$

$$best; y = (\cdot Ui)$$

# **Neural Tangent Kernel**

# Recipe for designing new kernels

$$f_{ ext{NN}}\left( heta_{ ext{NN}},x
ight) \gg k\left(x,x'
ight) = \mathbb{E}_{ heta_{ ext{NN}} \sim \mathcal{W}}\left[\left\langle \frac{\partial f_{ ext{NN}}\left( heta_{ ext{NN}},x
ight)}{\partial heta_{ ext{NN}}}, \frac{\partial f_{ ext{NN}}\left( heta_{ ext{NN}},x'
ight)}{\partial heta_{ ext{NN}}}
ight
angle
ight]$$

### Transform a neural network of any architecture to a kernel!

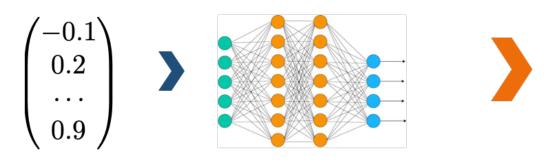
Fully-connected NN → Fully-connected NTK

Convolutional NN → Convolutional NTK

Graph NN → Graph NTK

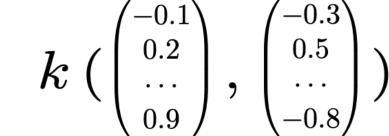
••••

# **Fully-Connect NTK**



**Features** 

**FC NN** 

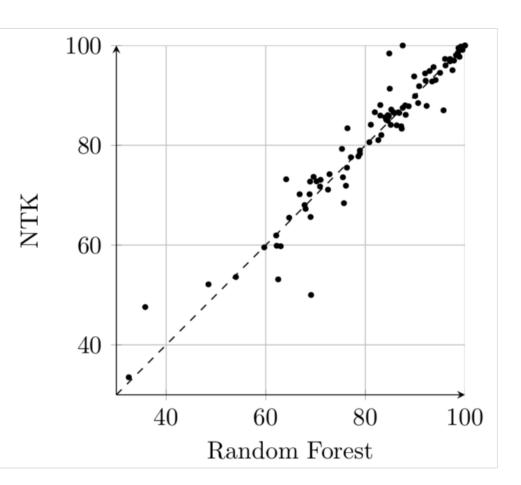


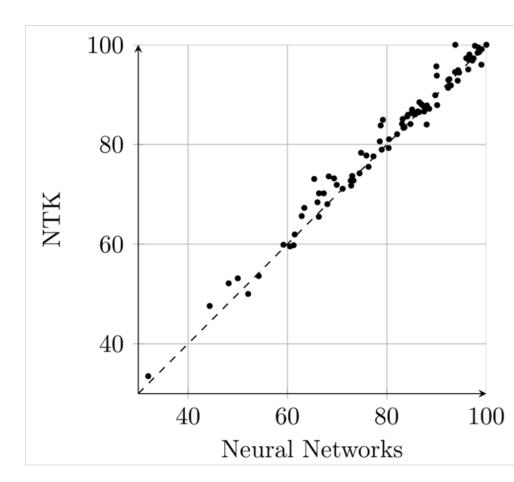
**FC NTK** 

		Avg Rank			
50 -		38			
40 -			35		
30 -	28	-	33		
20 -	_	_	_	_	
10 -					
	FC NTK FC NN		Random Forest	RBF Kernel	

Classifier	Avg Acc	P95	РМА
FC NTK	82%	<b>72%</b>	96%
FC NN	81%	60%	95%
Random Forest	82%	68%	95%
RBF Kernel	81%	<b>72</b> %	94%

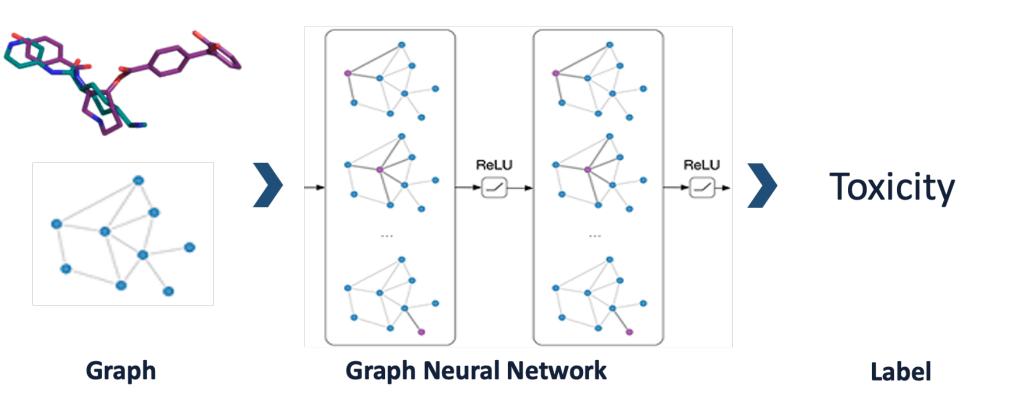
# **Pairwise Comparisons**





Classification Accuracy

# **Graph Neural Network**



# **Graph Neural Tangent Kernel**



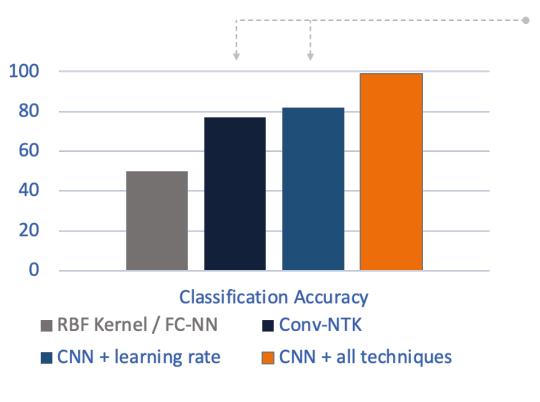
Graph Graph NN

**Graph NTK** 

	Method	COLLAB	IMDB-B	IMDB-M	PTC
GNN	GCN	79%	74%	51%	64%
	GIN	80%	75%	52%	65%
GK	WL	79%	74%	51%	60%
	GNTK	84%	77%	53%	68%

# What are left open?

### **CIFAR-10 Image Classification**



### **Open Problems:**

#### Why there is a gap:

finite-width? learning rate?

#### **Understanding techniques:**

batch-norm dropout data-augmentation

...

# Deep Learning Generalization



### **Measure of Generalization**

Generalization: difference in performance on train vs. test.

$$\frac{1}{n} \sum_{i=1}^{n} \mathcal{E}(f(x_i), y_i) - \mathbb{E}_{(x,y) \sim \mathcal{D}}[\mathcal{E}(f(x), y)]$$

Assumption  $(x_i, y_i)$   $i.i.d. \sim \mathcal{D}$ 

### Problems with the theoretical idealization

Data is not identically distributed:

- Images (Imagenet) are scraped in slightly different ways
- Data has systematic bias (e.g., patients are tested based on symptoms they exhibit)
- Data is result of interaction (reinforcement learning)
- Domain / distribution shift

### **Meta Theorem of Generalization**

**Meta theorem of generalization:** with probability  $1 - \delta$  over the choice of a training set of size n, we have

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

109 ([F])

# Some measures of complexity:

- (Log) number of elements
- VC (Vapnik-Chervonenkis) dimension
- Rademacher complexity
- PAC-Bayes
- ...

# Classical view of generalization

**Decoupled** view of generalization and optimization:

- Optimization: find a global minimum:  $\min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^{m} \ell(f(x_i), y_i)$
- Generalization: how well does the global optimizer generalize

**Practical implications:** to have a good generalization, make sure  $\mathcal{F}$  is not too "complex".

### Strategies:

- **Direct capacity control:** bound the size of the network / amount of connections, clip the weights, etc.
- Regularization: add a penalty term for "complex" predictors: weight decay ( $\ell_2$  norm), dropout, etc.

# Techniques for Improving Generalization



# **Weight Decay**

$$\cos f = \frac{\lambda}{2} 1104^2$$

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

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

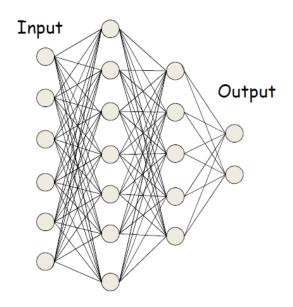
Adam M Mant uniform shrinkase for all  $G_{i}^{(t+1)} \leq (I - G_{i}(t)\lambda) G_{i}(t)$   $= G_{i}^{(t+1)} \vee f_{i}^{(t+1)} \vee f_{i}^{(t+1)}$ 

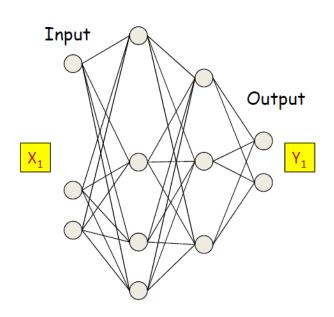
# **Dropout**

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.



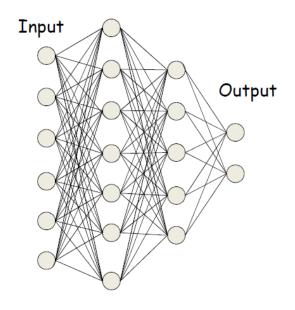


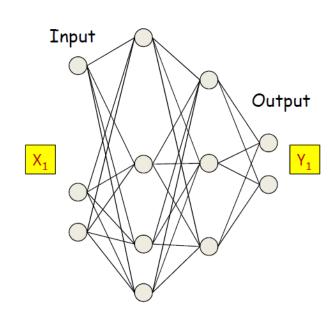
# **Dropout**

Dropout changes the scale of the output neuron:

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

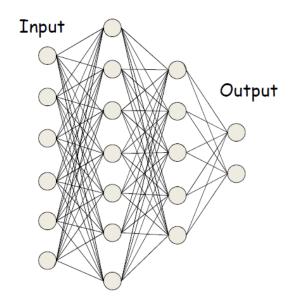
**Test time:**  $y = \alpha \sigma(Wx)$  to match the scale

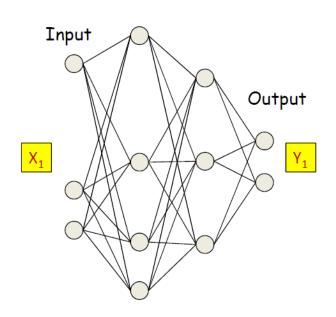




# **Understanding Dropout**

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

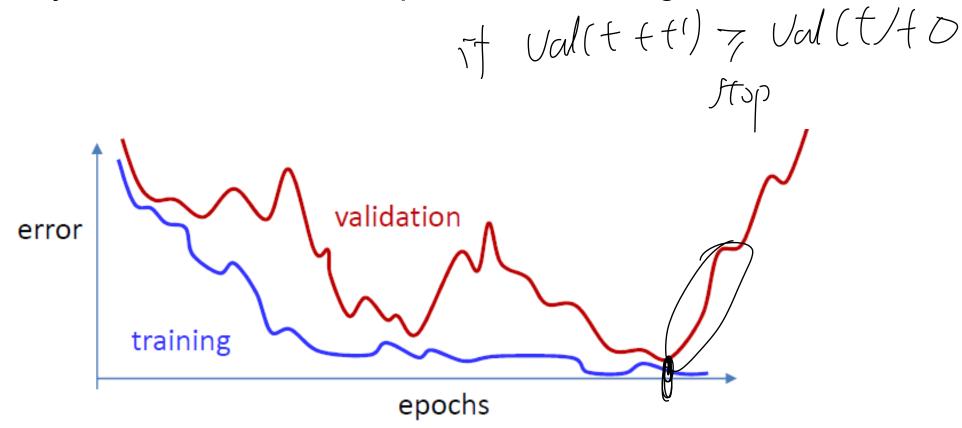




# **Early Stopping**

Vien A steatons as a hyperparameter

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



# **Data Augmentation**

dat a - centui

### Depend on data types.

JIN

### 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

• 
$$\hat{x} = \lambda x_i + (1 - \lambda)x_j$$

• 
$$\hat{y} = \lambda y_i + (1 - \lambda)y_i$$

•  $\lambda \sim \text{Beta}(0.2)$ 

# **Data Augmentation**

### Depend on data types.

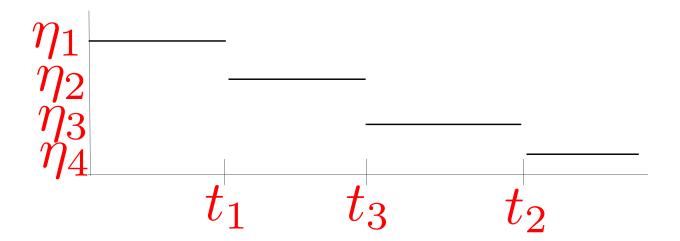
### Natural language processing:

- 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 schedule



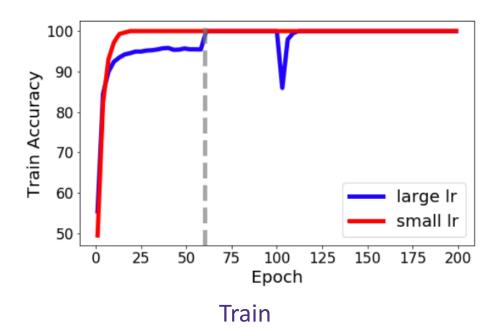
warm, cos, small

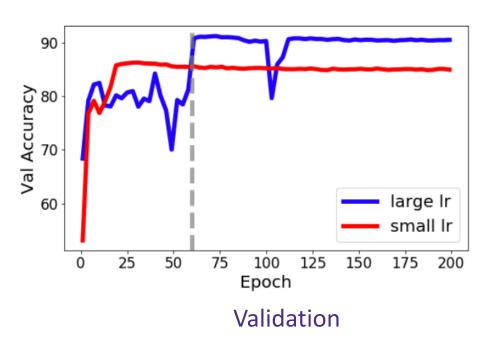
# Learning rate scheduling

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)

• . . .