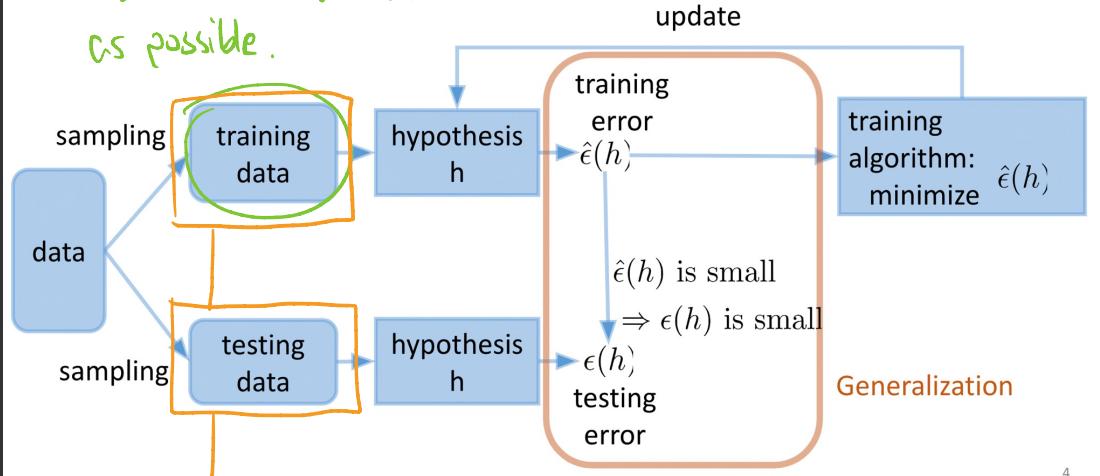


## \* Learning Theory

### Learning Theory

Required all situations  
as possible.

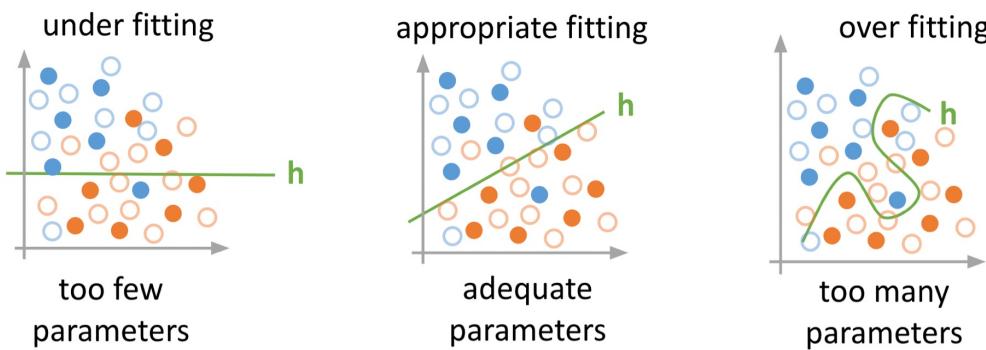


## \* Overfitting

### Learning Theory

- Over fitting :  $\hat{\epsilon}(h)$  small, but  $\epsilon(h)$  large
- What cause over-fitting?

● ● training data  
○ ○ testing data



$$\Sigma(h) \leq \hat{\Sigma}(h) + \sqrt{\frac{8}{n} \log\left(\frac{4(2n)^d}{\delta}\right)}$$

where  $d$  is VC Dimension (model. complexity),  $n$  : num. of training instances.

$n, d$  要要平衡

## \* VC Dimension.

- 1  $O(W)$ ,  $W$ : num. of params. for linear model
  - 2  $O(LW \cdot \log W)$ ,  $L$ : num. of layer,  $W$ : num. of params
- full-connected neural network.

## • Growth Function

$$\tau(x_1, \dots, x_n) = \left| \left\{ (h(x_1), \dots, h(x_n)) \mid h \in H \right\} \right|$$

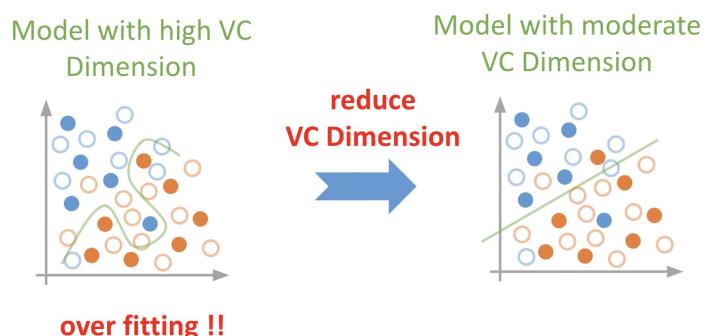
$$\tau_H(n) = \max_{(x_1, \dots, x_n)} \tau(x_1, \dots, x_n)$$

data samples      hypothesis set

## • VC Dimension

$$d(H) = \max\{n : \tau_H(n) = 2^n\}$$

## \* VC Bound.



## \* Generalization

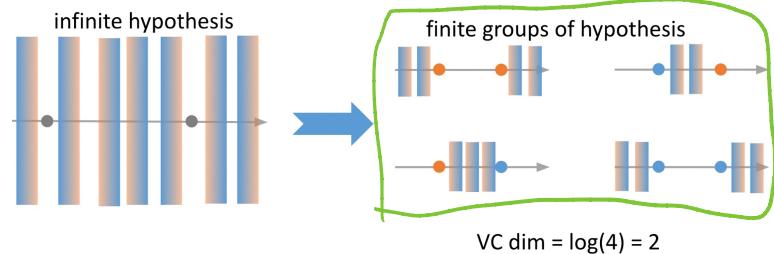
Better situation  $\hat{\epsilon}(h) \sim \epsilon(h)$

## VC Dimension

## Model Complexity

Categorized the hypothesis set into finite groups

- ex: 1D linear model

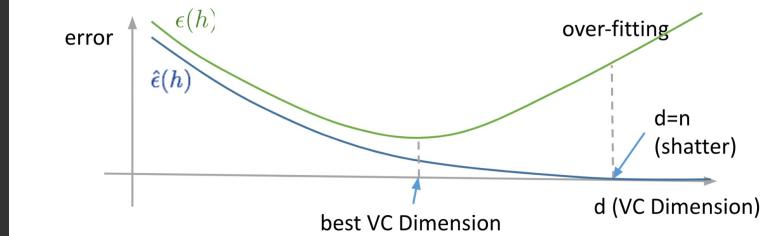


Ex  $\tau(x_1, x_2) = |\{(h(x_1), h(x_2)) \mid h \in H\}| = 4$   
 $\{(h(x_1), h(x_2))\} = \{(1,1), (0,1), (1,0), (0,0)\} \rightarrow$   
 ⇒ 可以有幾種狀況的維度

## VC Bound

- For a given dataset ( $n$  is constant), search for the best VC Dimension

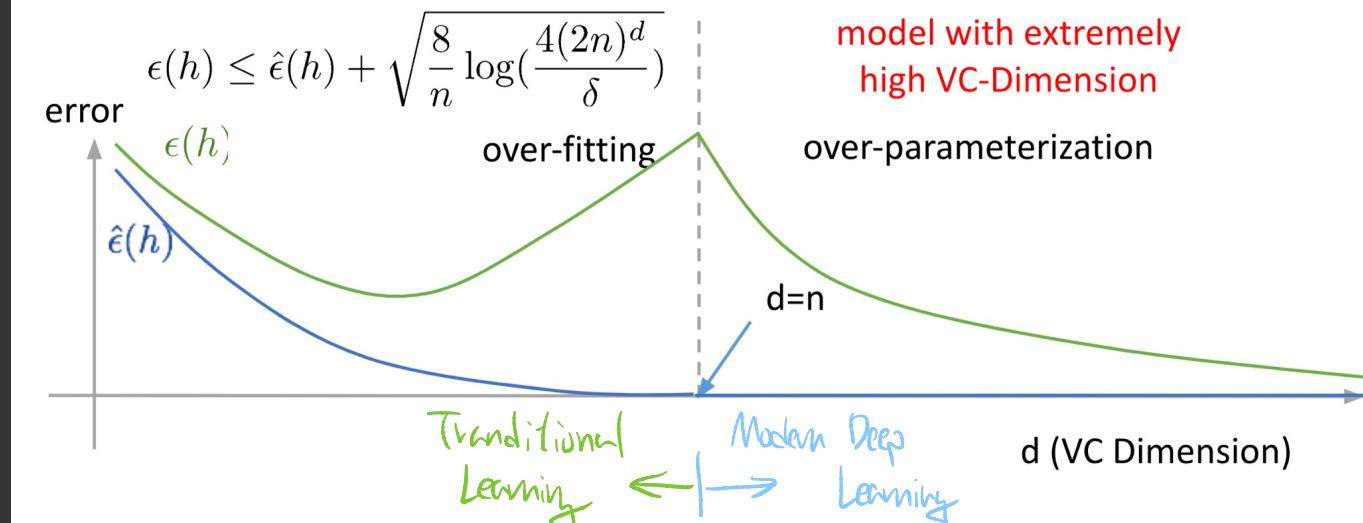
$$\epsilon(h) \leq \hat{\epsilon}(h) + \sqrt{\frac{8}{n} \log(\frac{4(2n)^d}{\delta})}$$



$\hat{\epsilon}(h)$ : training error ;  $\epsilon(h)$ : testing error

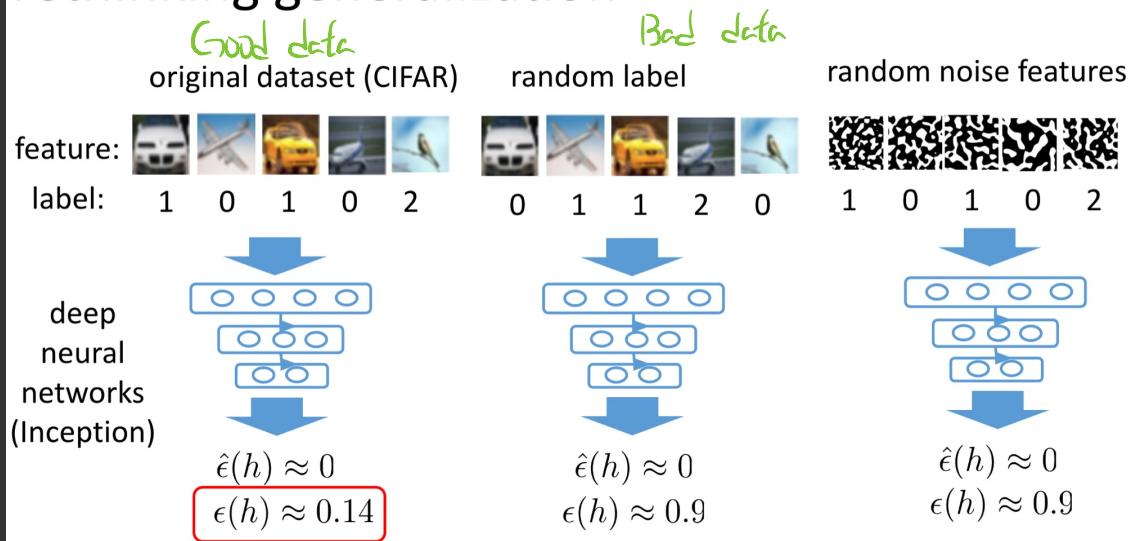
\* Modern machine learning

## Reconciling modern machine learning and the bias-variance trade-off



Fur modern model due to that  
hardware can handle very large data now.

## Understanding deep learning requires rethinking generalization



基本上有資料都能練成一個模型，但資料狀況不佳會導致 testing 畏到  
人工智慧是人工智慧

\* Rethinking generalization for modern model

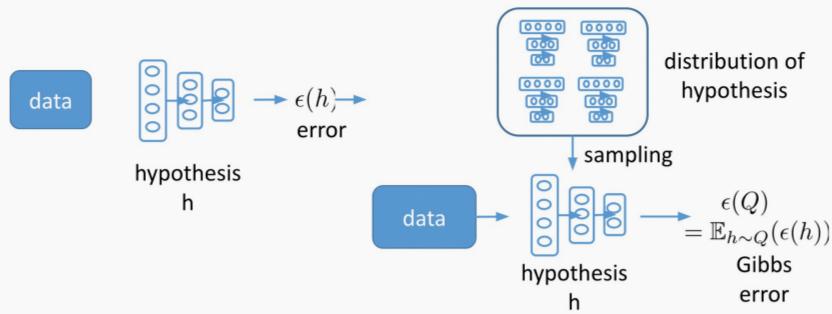
Regularization 對於不同模型效果不一，但改善  $\hat{\epsilon}(h)$  效果不大

## \* Pac - Bayesian Framework.

Talking about how regularization affect models.

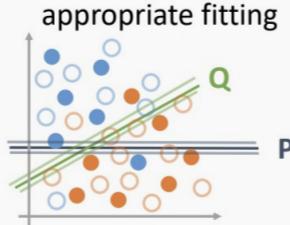
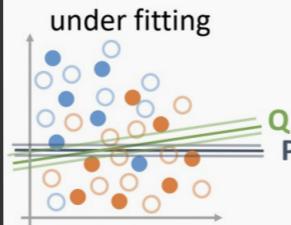
### PAC-Bayesian Bound

- Deterministic Model
- Stochastic Model (Gibbs Classifier)

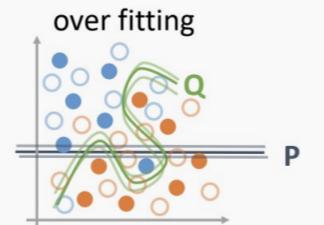


### PAC-Bayesian Bound

$$\epsilon(Q) \leq \hat{\epsilon}(Q) + \sqrt{\frac{KL(Q||P) + \log(\frac{n}{\delta}) + 2}{2n - 1}}$$



- ● training data
- ○ testing data



Stop training in a lower  
 $|\epsilon(Q) - \hat{\epsilon}(Q)|$

### Experiments

original M-NIST

1	2	3	7	9
0	0	0	1	1

random label

1	2	3	7	9
1	0	1	0	1

Training	0.028
Testing	0.034
VC Bound	26m
Pac-Bayesian Bound	0.161

0.112
0.503
26m
1.352

Pac - Bayesian Bound 相對適合 Deep Learning.

## \* Comparison

# How to Overcome Over-fitting?

$$\epsilon(h) \leq \hat{\epsilon}(h) + \sqrt{\frac{8}{n} \log\left(\frac{4(2n)^d}{\delta}\right)}$$

## Traditional Machine Learning

- reduce the number of parameters
- weight decay
- early stop
- data augmentation

資料僵化

$$KL(\hat{\epsilon}(Q) \| \epsilon(Q)) \leq \frac{KL(Q \| P) + \log\left(\frac{n}{\delta}\right)}{n-1}$$

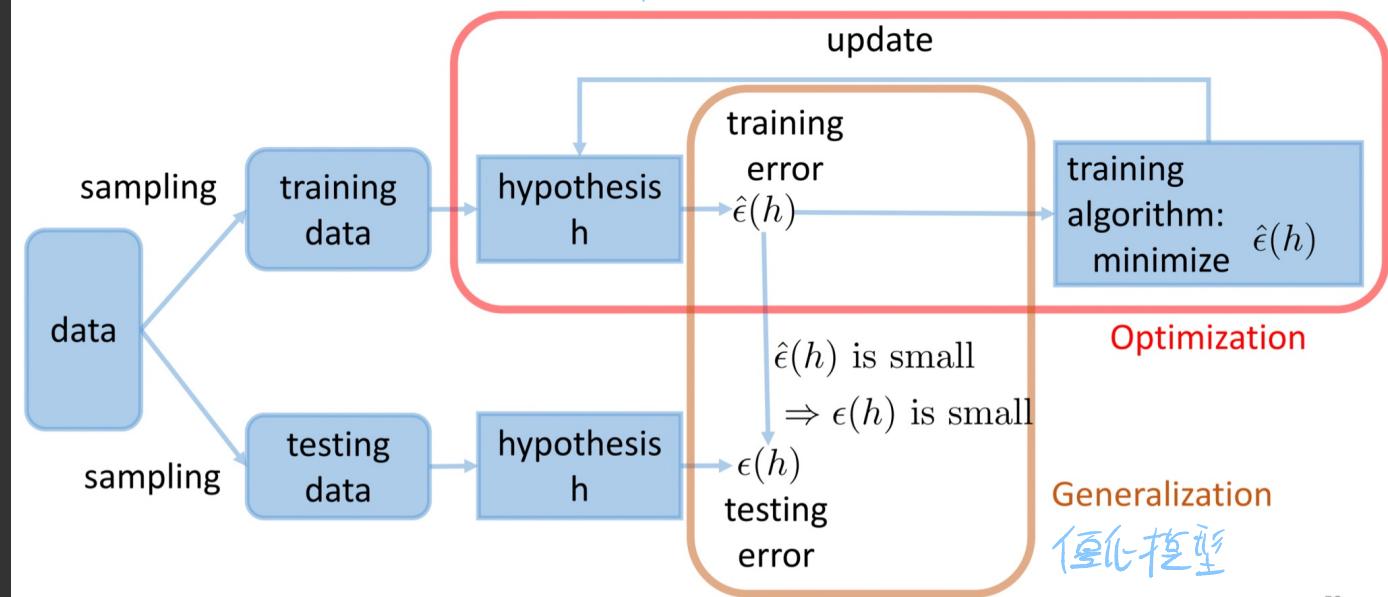
## Modern Deep Learning *Most Different*

- ~~reduce the number of parameters~~
- weight decay ?
- early stop ?
- data augmentation ?
- ~~improving data quality~~
- starting from a good P

## • Neural Tangent kernel

# Learning Theory

僵化訓練效能



# Analysis of Generalization

- Theorem 5.1. v.s. VC Bound v.s. PAC-Bayesian Bound

- Theorem 5.1.

$$L_{\mathcal{D}}(f_{\mathbf{W}(k), \mathbf{a}}) \leq \sqrt{\frac{2\mathbf{y}^T(\mathbf{H}^\infty)^{-1}\mathbf{y}}{n}} + O\left(\sqrt{\frac{\log \frac{n}{\lambda_0 \delta}}{n}}\right)$$

- Only depends on training data
- Model doesn't need to be trained
- Can only be applied to over-parameterized 2-layer ReLU NN

- VC Bound

$$\epsilon(h) \leq \hat{\epsilon}(h) + \sqrt{\frac{8}{n} \log\left(\frac{4(2n)^d}{\delta}\right)}$$

- Only depends on model
- Can not be applied to over-parameterization NN

- PAC-Bayesian Bound

$$\epsilon(Q) \leq \hat{\epsilon}(Q) + \sqrt{\frac{KL(Q||P) + \log(\frac{n}{\delta}) + 2}{2n - 1}}$$

- Depends on both model and training data
- Model needs to be trained