

Outline

- Feedforward Neural Networks
- Deep learning in the past 50th years
- Some popular neural networks

Feedforward Neural Networks

- ▶ Definition of Neural Networks
 - Forward propagation
 - Types of units
 - Capacity of neural networks
- ▶ How to train neural nets:
 - Loss function
 - Backpropagation with gradient descent
- ▶ More recent techniques:
 - Dropout
 - Batch normalization
 - Unsupervised Pre-training

A diagram of a feedforward neural network. It shows four layers of nodes. The input layer has nodes labeled $x_1, \dots, x_j, \dots, x_d$. Above it is a hidden layer with three nodes. Above that is another hidden layer with two nodes. The output layer consists of two nodes labeled $f(x)$ and \dots . Arrows indicate the connections between nodes in adjacent layers, showing how information flows from left to right through the network.

Artificial Neuron

- Neuron pre-activation (or input activation):

$$a(\mathbf{x}) = b + \sum_i w_i x_i = b + \mathbf{w}^\top \mathbf{x}$$

- Neuron output activation:

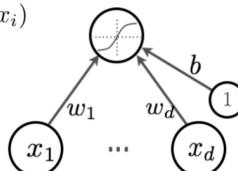
$$h(\mathbf{x}) = g(a(\mathbf{x})) = g(b + \sum_i w_i x_i)$$

where

\mathbf{w} are the weights (parameters)

b is the bias term

$g(\cdot)$ is called the activation function

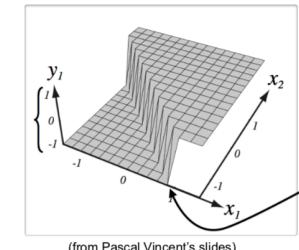


Artificial Neuron

- Output activation of the neuron:

$$h(\mathbf{x}) = g(a(\mathbf{x})) = g(b + \sum_i w_i x_i)$$

Range is determined by $g(\cdot)$



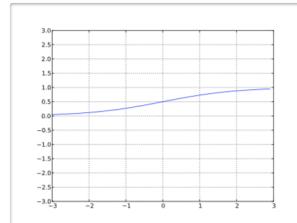
Bias only changes the position of the riff

Activation Function

- Sigmoid activation function:

- Squashes the neuron's output between 0 and 1

$$g(a) = \text{sigm}(a) = \frac{1}{1+\exp(-a)}$$



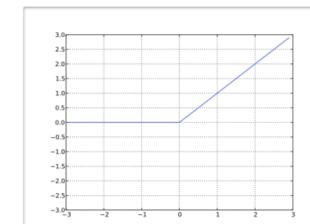
- Always positive
- Bounded
- Strictly Increasing

- Rectified linear (ReLU) activation function:

- Bounded below by 0 (always non-negative)
倾向于产生稀疏活动的单元

- Tends to produce units with sparse activities
- Not upper bounded
- Strictly increasing

$$g(a) = \text{relin}(a) = \max(0, a)$$



Single Hidden Layer Neural Net

- Hidden layer pre-activation:

$$\mathbf{a}(\mathbf{x}) = \mathbf{b}^{(1)} + \mathbf{W}^{(1)}\mathbf{x}$$

$$(a(\mathbf{x})_i = b_i^{(1)} + \sum_j W_{i,j}^{(1)} x_j)$$

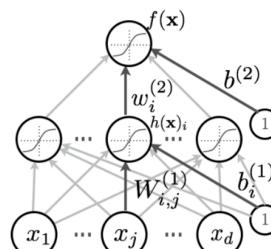
- Hidden layer activation:

$$\mathbf{h}(\mathbf{x}) = \mathbf{g}(\mathbf{a}(\mathbf{x}))$$

- Output layer activation:

$$f(\mathbf{x}) = o\left(b^{(2)} + \mathbf{w}^{(2)^\top} \mathbf{h}^{(1)} \mathbf{x}\right)$$

Output activation function



Multilayer Neural Net

- Consider a network with L hidden layers.

- layer pre-activation for k>0

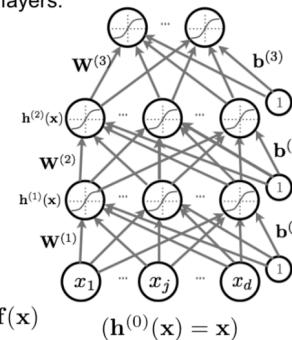
$$\mathbf{a}^{(k)}(\mathbf{x}) = \mathbf{b}^{(k)} + \mathbf{W}^{(k)}\mathbf{h}^{(k-1)}(\mathbf{x})$$

- hidden layer activation from 1 to L:

$$\mathbf{h}^{(k)}(\mathbf{x}) = \mathbf{g}(\mathbf{a}^{(k)}(\mathbf{x}))$$

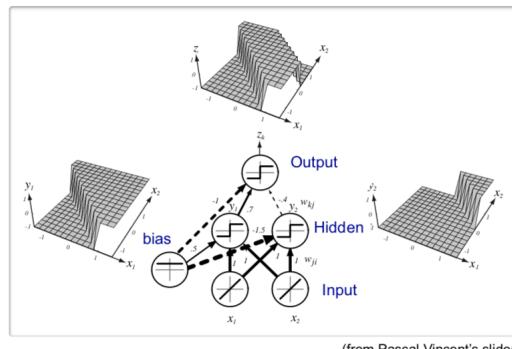
- output layer activation (k=L+1):

$$\mathbf{h}^{(L+1)}(\mathbf{x}) = \mathbf{o}(\mathbf{a}^{(L+1)}(\mathbf{x})) = \mathbf{f}(\mathbf{x})$$



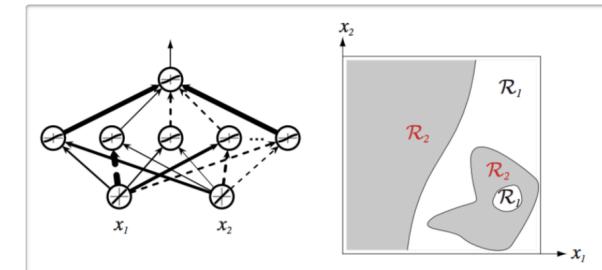
Capacity of Neural Nets

- Consider a single layer neural network



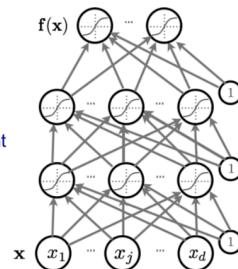
Capacity of Neural Nets

- Consider a single layer neural network



Feedforward Neural Networks

- ▶ How neural networks predict $f(x)$ given an input x :
 - Forward propagation
 - Types of units
 - Capacity of neural networks
 - ▶ How to train neural nets:
 - Loss function
 - Backpropagation with gradient descent
 - ▶ More recent techniques:
 - Dropout
 - Batch normalization
 - Unsupervised Pre-training



Training

- Empirical Risk Minimization:

$$\arg \min_{\theta} \frac{1}{T} \sum_t l(f(\mathbf{x}^{(t)}; \theta), y^{(t)}) + \lambda \Omega(\theta)$$

The diagram shows the expression above with two curly braces. The first brace, under the term $\sum_t l(f(\mathbf{x}^{(t)}; \theta), y^{(t)})$, is labeled "Loss function". The second brace, under the term $\lambda \Omega(\theta)$, is labeled "Regularizer".

- To train a neural net, we need

- **Loss function:** $l(f(\mathbf{x}^{(t)}; \boldsymbol{\theta}), y^{(t)})$
 - A procedure to **compute gradients**: $\nabla_{\boldsymbol{\theta}} l(f(\mathbf{x}^{(t)}; \boldsymbol{\theta}), y^{(t)})$
 - **Regularizer** and its gradient: $\Omega(\boldsymbol{\theta}), \nabla_{\boldsymbol{\theta}} \Omega(\boldsymbol{\theta})$

Stochastic Gradient Descend

- Perform updates after seeing each example

- Initialize: $\theta \equiv \{\mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \dots, \mathbf{W}^{(L+1)}, \mathbf{b}^{(L+1)}\}$

- For $t=1:T$

- for each training example $(\mathbf{x}^{(t)}, y^{(t)})$

$$\Delta = -\nabla_{\theta} l(f(\mathbf{x}^{(t)}; \theta), y^{(t)}) - \lambda \nabla_{\theta} \Omega(\theta)$$

$$\theta \leftarrow \theta + \alpha \Delta$$

Training epoch
= Iteration of all examples

- To train a neural net, we need:

- **Loss function:** $l(\mathbf{f}(\mathbf{x}^{(t)}; \boldsymbol{\theta}), y^{(t)})$
 - A procedure to **compute gradients**: $\nabla_{\boldsymbol{\theta}} l(\mathbf{f}(\mathbf{x}^{(t)}; \boldsymbol{\theta}), y^{(t)})$
 - **Regularizer** and its gradient: $\Omega(\boldsymbol{\theta}), \nabla_{\boldsymbol{\theta}} \Omega(\boldsymbol{\theta})$

Weight Decay 权重衰减

$$\arg \min_{\boldsymbol{\theta}} \frac{1}{T} \sum_t l(f(\mathbf{x}^{(t)}; \boldsymbol{\theta}), y^{(t)}) + \lambda \Omega(\boldsymbol{\theta})$$

- L2 regularization

$$\Omega(\theta) = \sum_k \sum_i \sum_j \left(W_{i,j}^{(k)} \right)^2 = \sum_k \| \mathbf{W}^{(k)} \|_F^2$$

- L1 regularization

$$\Omega(\theta) = \sum_k \sum_i \sum_j |W_{i,j}^{(k)}|$$

- Only applies to weights, not biases (weight decay)

Model Selection

- Training Protocol:
 - Train your model on the **Training Set** $\mathcal{D}^{\text{train}}$
 - For model selection, use **Validation Set** $\mathcal{D}^{\text{valid}}$
 - Hyper-parameter search: hidden layer size, learning rate, number of iterations/epochs, etc.
 - Estimate generalization performance using the **Test Set** $\mathcal{D}^{\text{test}}$
- Generalization is the behavior of the model on **unseen examples**. **泛化**

Early Stopping

- To select the number of epochs, stop training when validation set error increases (with some look ahead).



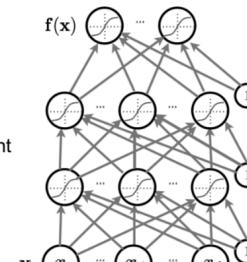
Mini-batch, Momentum

- Make updates based on a mini-batch of examples (instead of a single example):
 - the gradient is the average regularized loss for that mini-batch
 - can give a more accurate estimate of the gradient
 - can leverage matrix/matrix operations, which are more efficient
- **Momentum:** Can use an exponential average of previous gradients:

$$\bar{\nabla}_{\theta}^{(t)} = \nabla_{\theta} l(\mathbf{f}(\mathbf{x}^{(t)}), y^{(t)}) + \beta \bar{\nabla}_{\theta}^{(t-1)}$$
 - can get pass plateaus more quickly, by "gaining momentum"
 - 通过增加动量，可以快速越过高原

Feedforward Neural Networks

- How neural networks predict $f(x)$ given an input x :
 - Forward propagation
 - Types of units
 - Capacity of neural networks
- How to train neural nets:
 - Loss function
 - Backpropagation with gradient descent
- More recent techniques:
 - Dropout
 - Batch normalization



Best Practice

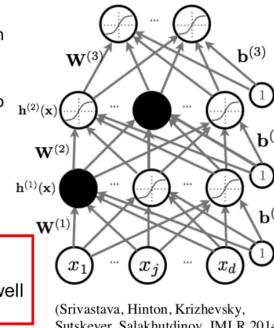
- Given a dataset D, pick a model so that:
 - You can achieve 0 training error—Overfit on the training set
- Regularize the model (e.g. using Dropout).
- SGD with momentum, batch-normalization, and dropout usually works very well.

Dropout

- Key idea:** Cripple neural network by removing hidden units stochastically

- each hidden unit is set to 0 with probability 0.5
- hidden units cannot co-adapt to other units
- hidden units must be more generally useful

- Could use a different dropout probability, but 0.5 usually works well



Dropout

- Use random binary masks $m^{(k)}$

layer pre-activation for $k > 0$

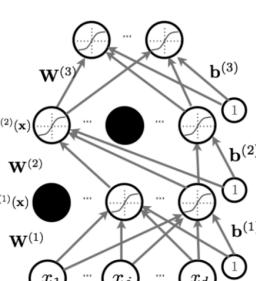
$$\mathbf{a}^{(k)}(\mathbf{x}) = \mathbf{b}^{(k)} + \mathbf{W}^{(k)}\mathbf{h}^{(k-1)}(\mathbf{x})$$

hidden layer activation ($k=1$ to L):

$$\mathbf{h}^{(k)}(\mathbf{x}) = g(\mathbf{a}^{(k)}(\mathbf{x})) \odot m^{(k)}$$

Output activation ($k=L+1$)

$$\mathbf{h}^{(L+1)}(\mathbf{x}) = o(\mathbf{a}^{(L+1)}(\mathbf{x})) = f(\mathbf{x})$$



(Srivastava, Hinton, Krizhevsky, Sutskever, Salakhutdinov, JMLR 2014)

Dropout at Test Time

- At test time, we replace the masks by their expectation 测试时使用mask的期望值

This is simply the constant vector 0.5 if dropout probability is 0.5

For single hidden layer: equivalent to taking the geometric average of all neural networks, with all possible binary masks

- Can be combined with unsupervised pre-training

- Beats regular backpropagation on many datasets

- Ensemble:** Can be viewed as a geometric average of exponential number of networks. 整合：可以看作指数个网络的几何平均

Batch Normalization

加速训练

- Normalizing the inputs will **speed up training** (Lecun et al. 1998)
 - could normalization be useful at the level of the hidden layers?
- Batch normalization is an attempt to do that (Ioffe and Szegedy, 2015)
 - each unit's pre-activation is normalized (mean subtraction, stddev division)
 - during training, mean and stddev is computed for each minibatch
 - backpropagation takes into account the normalization
 - at test time, the global mean / stddev is used

□ (Ioffe and Szegedy, ICML 2015)

Batch Normalization

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_1 \dots m\}$;

Parameters to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$

Learned linear transformation to adapt to non-linear activation function (γ and β are trained)

(Ioffe and Szegedy, ICML 2015)

Batch Normalization

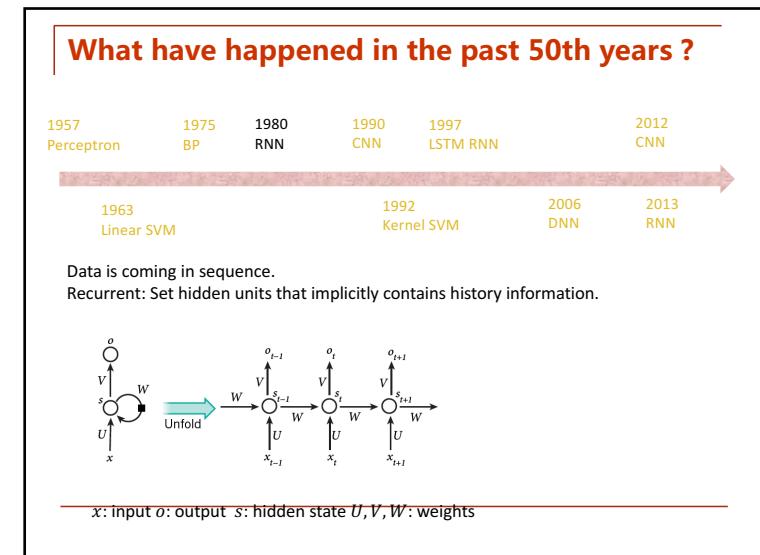
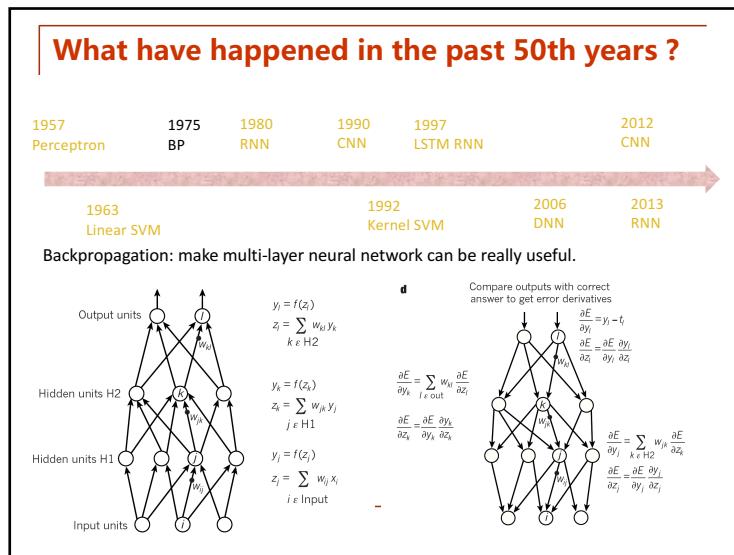
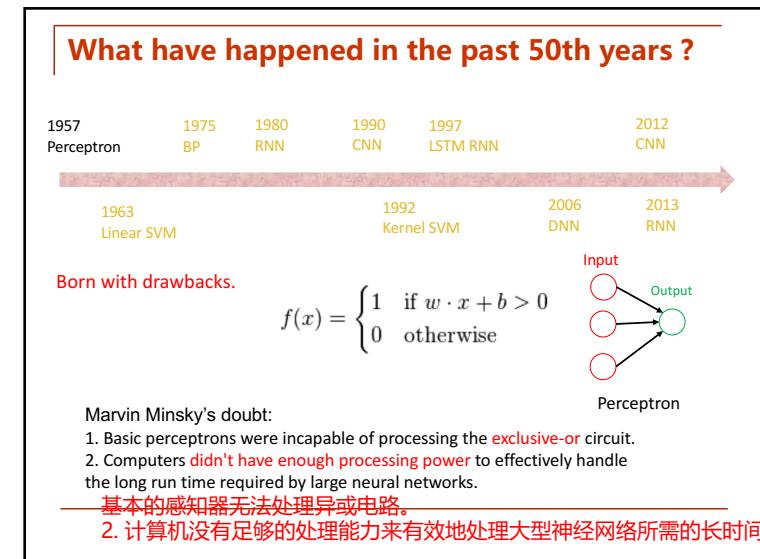
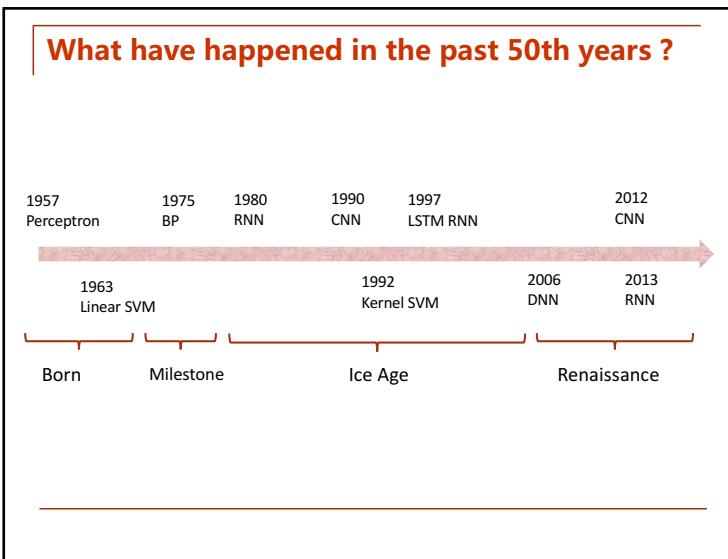


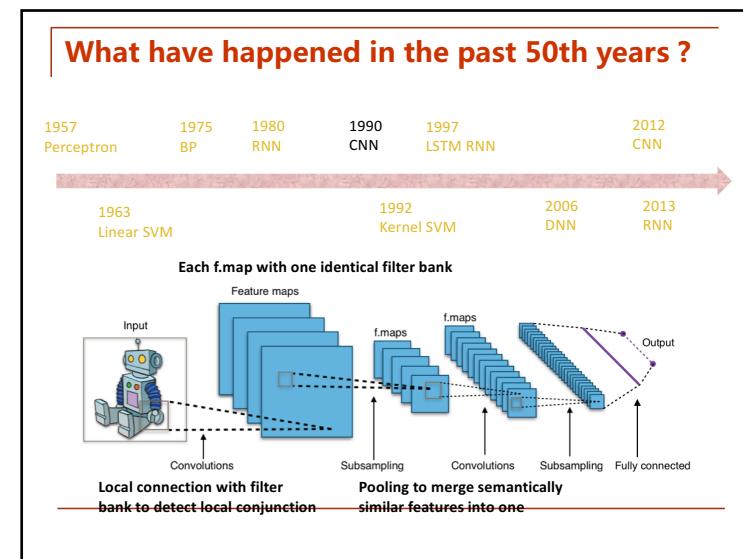
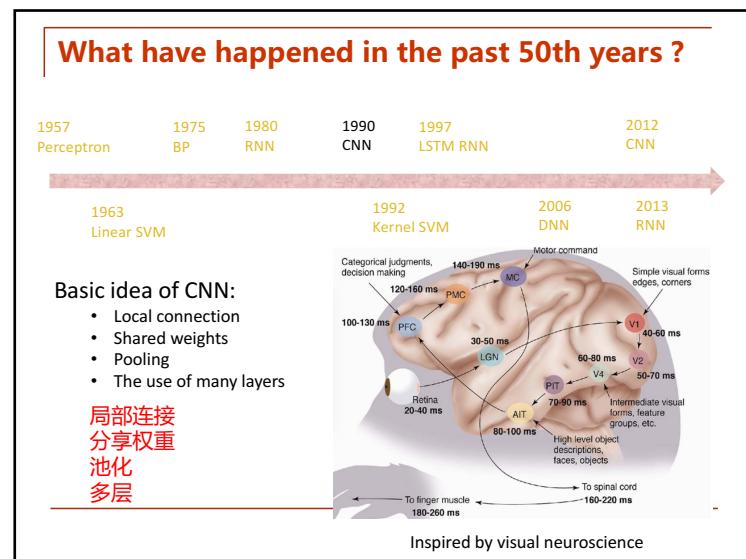
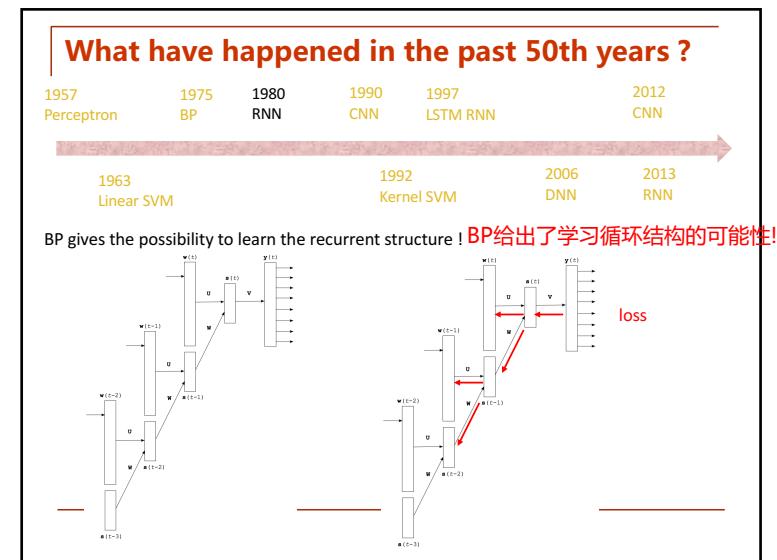
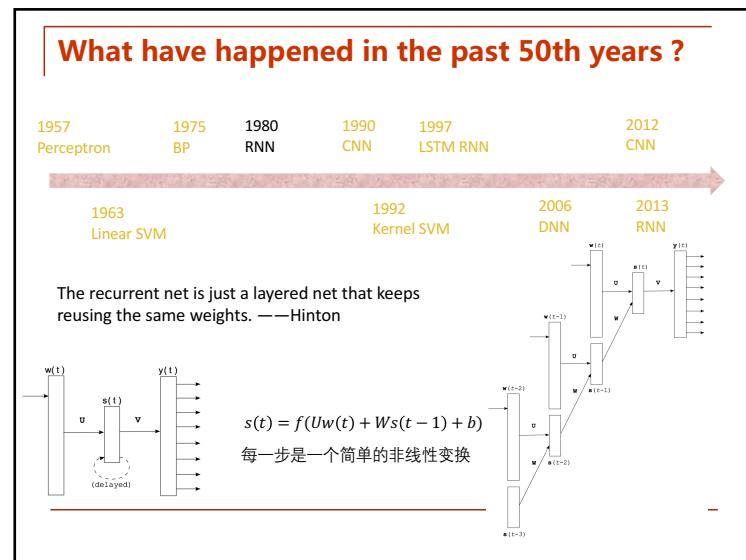
- Why normalize the pre-activation?
 - can help keep the pre-activation in a non-saturating regime (though the linear transform $y_i \leftarrow \gamma \hat{x}_i + \beta$ could cancel this effect)
- Use the **global mean and stddev** at test time.
 - removes the stochasticity of the mean and stddev **去除随机性**
 - requires a final phase where, from the first to the last hidden layer
 - propagate all training data to that layer
 - compute and store the global mean and stddev of each unit
 - for early stopping, could use a running average

(Ioffe and Szegedy, ICML 2015)

Outline

- Deep Neural Networks
- Deep learning in the past 50th years





What have happened in the past 50th years ?

Timeline:

- 1957 Perceptron
- 1975 BP
- 1980 RNN
- 1990 CNN
- 1997 LSTM RNN
- 2012 CNN

1963 Linear SVM

1992 Kernel SVM

$k(\mathbf{x}, \mathbf{x}') = \langle \varphi(\mathbf{x}), \varphi(\mathbf{x}') \rangle_V$

Kernel trick brings nonlinear
But mystery nonlinear...

Primal

$$\text{minimize } \frac{1}{n} \sum_{i=1}^n \zeta_i + \frac{\lambda}{2} \|w\|^2$$

Dual

$$\text{maximize } f(c_1 \dots c_n) = \sum_{i=1}^n c_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i c_i (x_i \cdot x_j) y_j c_j,$$

subject to $y_i(x_i \cdot w + b) \geq 1 - \zeta_i$ and $\zeta_i \geq 0$, for all i .

subject to $\sum_{i=1}^n c_i y_i = 0$, and $0 \leq c_i \leq \frac{1}{2n\lambda}$ for all i .

What have happened in the past 50th years ?

Timeline:

- 1957 Perceptron
- 1975 BP
- 1980 RNN
- 1990 CNN
- 1997 LSTM RNN
- 2012 CNN

1963 Linear SVM

1992 Kernel SVM

2006 DNN

2013 RNN

Dark days for neural networks.
Yes, we have a powerful model;
but, we do not have a useful
learning algorithm !

Unknown target function $f: \mathcal{X} \rightarrow \mathcal{Y}$

Training examples $D: (x_1, y_1), \dots, (x_N, y_N)$

Learning Algorithm A

Final hypothesis $g \approx f$

Supervised: Known all y_n
Unsupervised: Known no y_n
Semi-supervised: Known some y_n

• BP do not work well for deep architecture
• Too many computation costs

Hypothesis Set \mathcal{H}

How to optimize?

What is the model of the data?

What have happened in the past 50th years ?

Timeline:

- 1957 Perceptron
- 1975 BP
- 1980 RNN
- 1990 CNN
- 1997 LSTM RNN
- 2012 CNN

1963 Linear SVM

1992 Kernel SVM

2006 DNN

2013 RNN

Here comes renaissance.

Big data

Computing Power

Enough training dataset

Complex optimization could be possible

What have happened in the past 50th years ?

Timeline:

- 1957 Perceptron
- 1975 BP
- 1980 RNN
- 1990 CI
- 1997 Le Cun, Hinton, Bengio, Montréal
- 2006 DNN
- 2012 CNN
- 2013 RNN

1963 Linear SVM

Renaissance of Deep Neural Network 深度神经网络的复兴

Unsupervised pre-training + fine tune

features

input

reconstruction of input

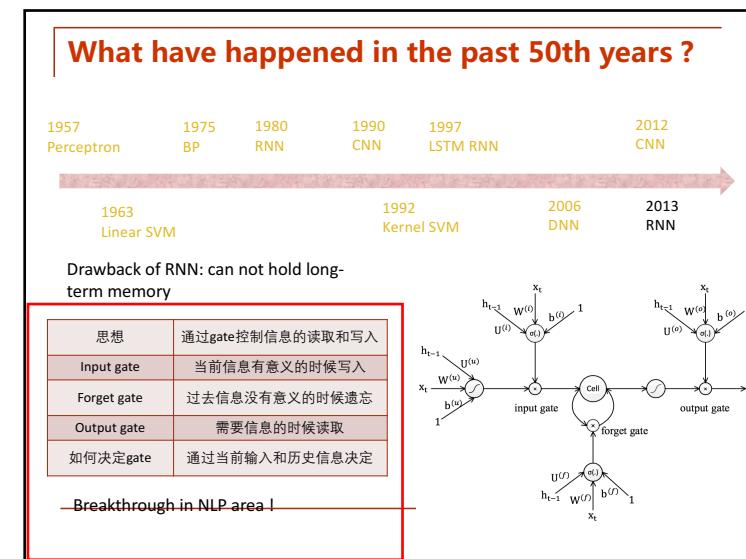
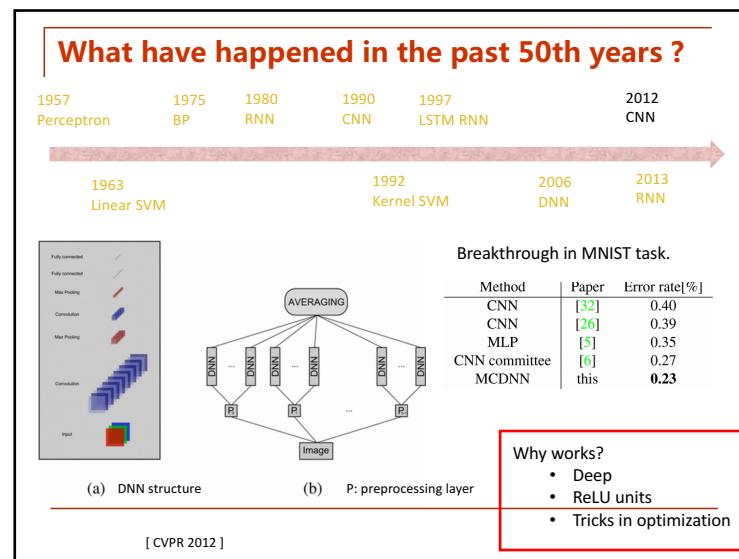
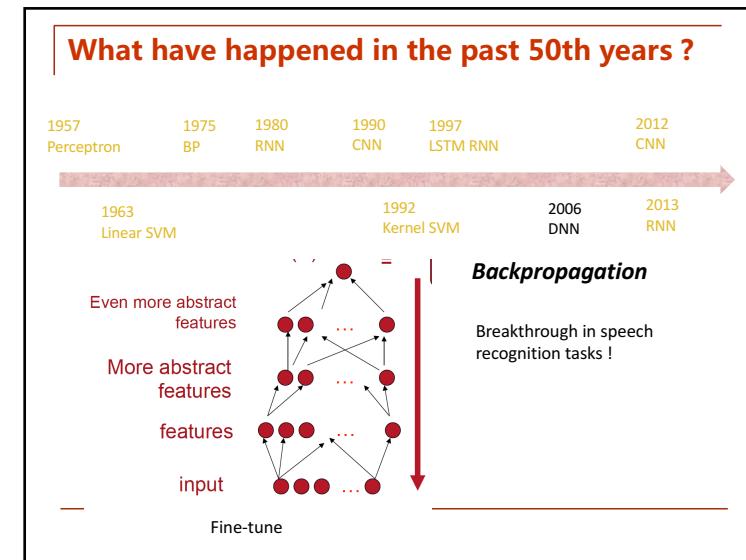
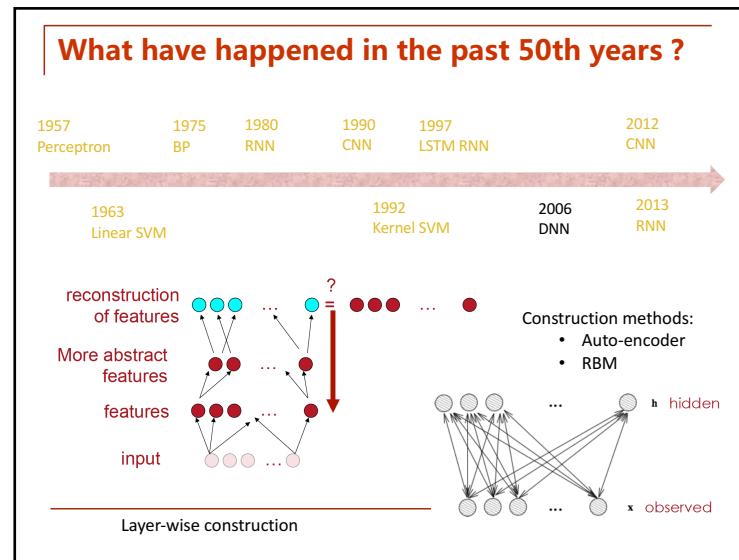
features

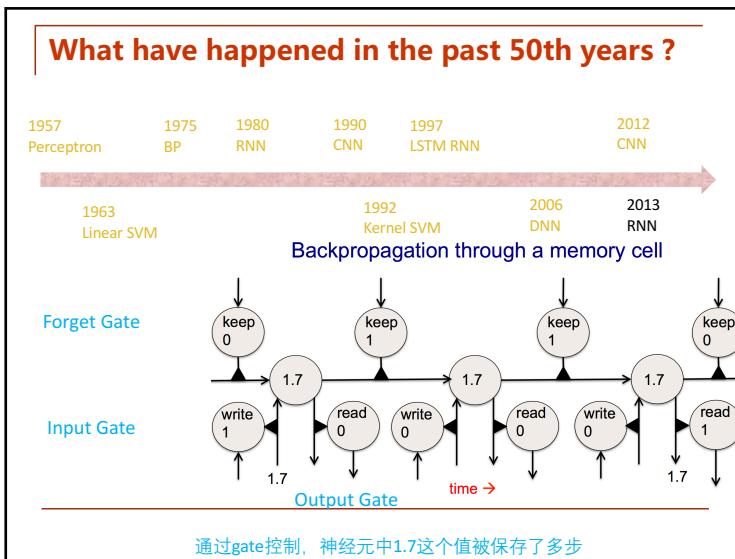
input

?

=

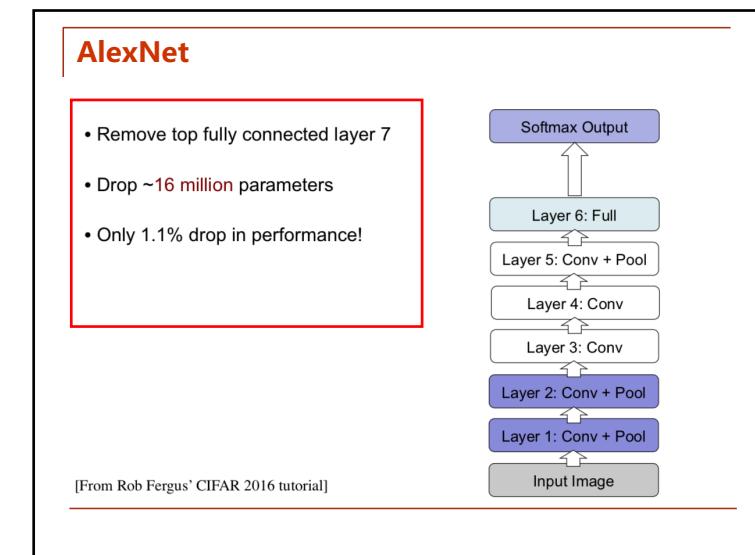
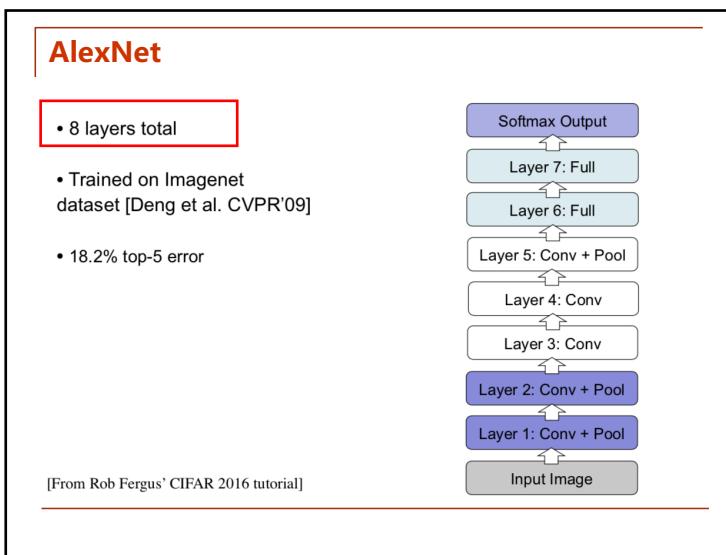
Construct a layer by supervised method. (Auto-encoder here)

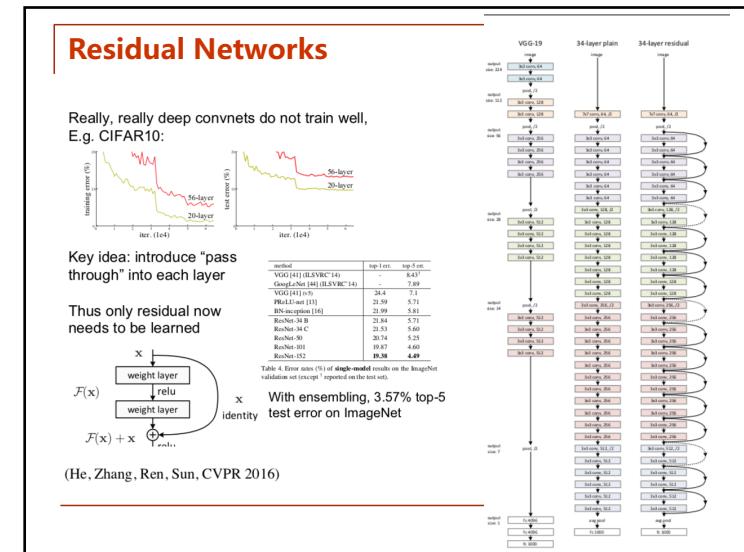
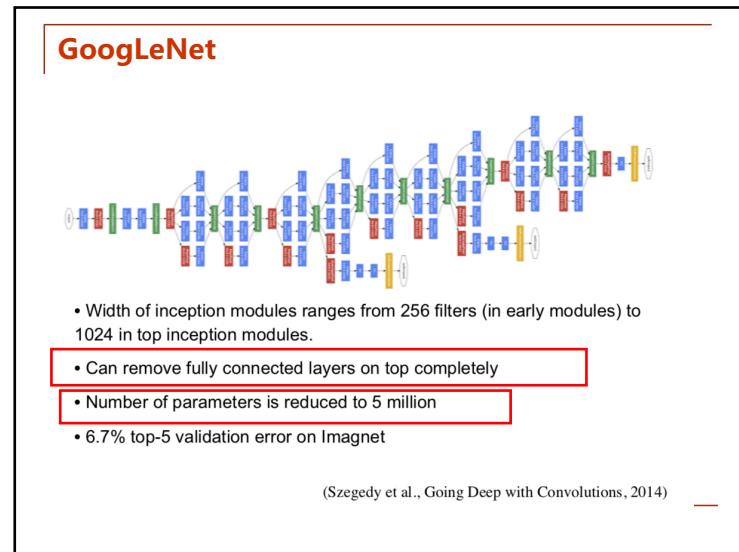
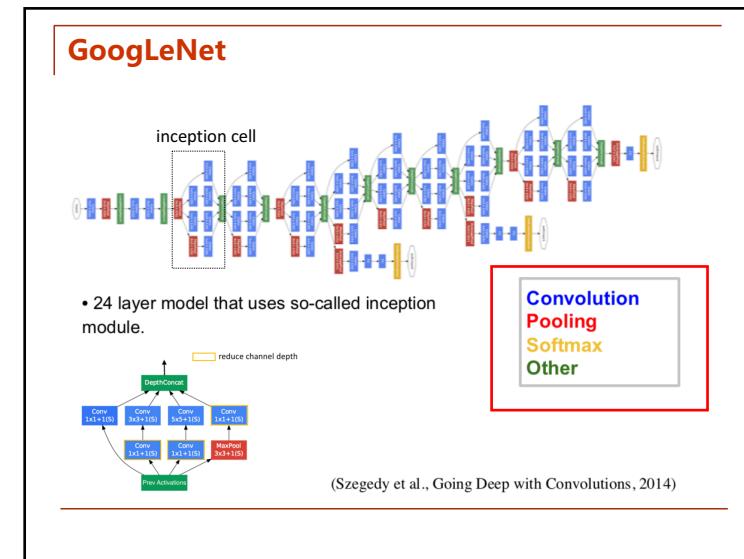
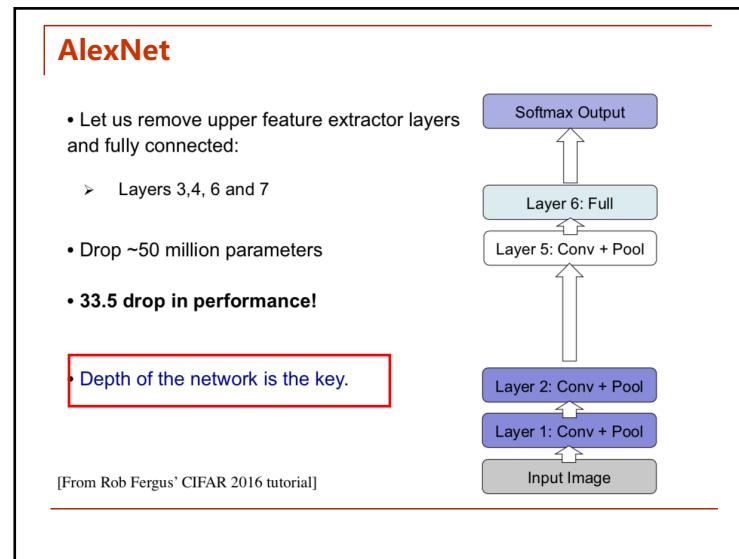




Outline

- Feedforward Neural Networks
- Deep learning in the past 50th years
- Some popular neural networks





Future?

- DL has big break through in visual, speech, NLP areas

那么关系数据呢

- How about the relational data ?

- Network Embedding 网络嵌入
- Recommendation 推荐
- Heterogeneous network matching 异构网络匹配

Questions?