

CSE 152: Computer Vision

Hao Su

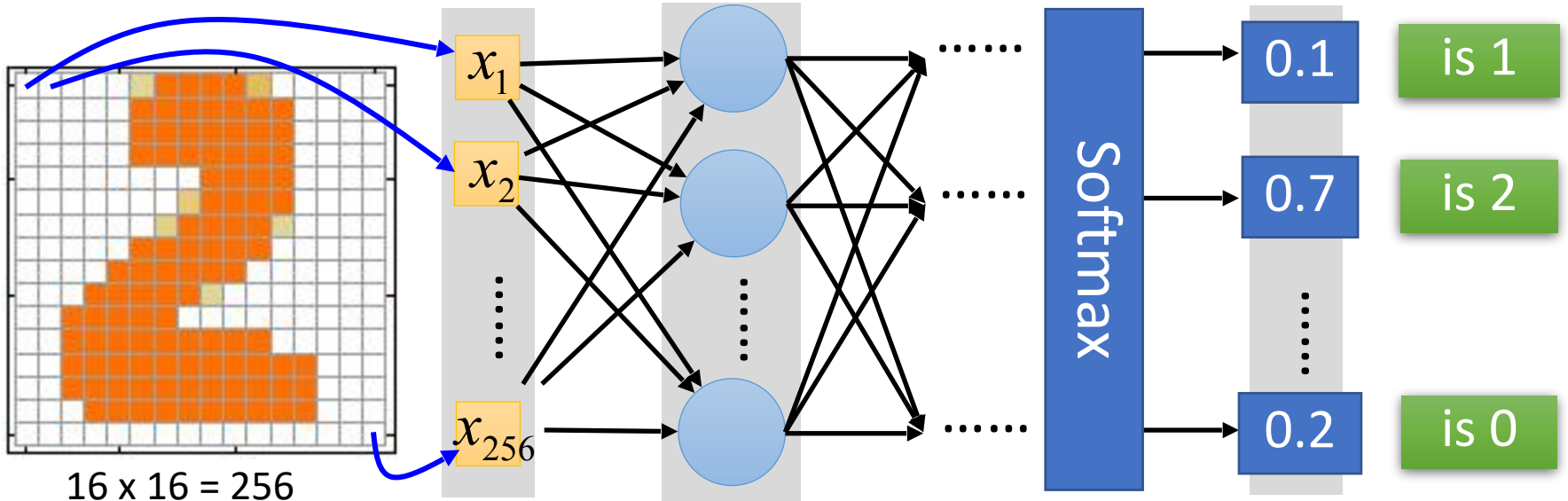
Lecture 8: Statistical and Optimization Perspectives of Deep Learning



Optimization of Neural Network

How to set network parameters

$$\theta = \{W_1, b_1, \dots, W_n, b_n\}$$



16 x 16 = 256

Ink \rightarrow 1

No ink \rightarrow 0

Set the network parameters θ such that

Input x_1 has the maximum value

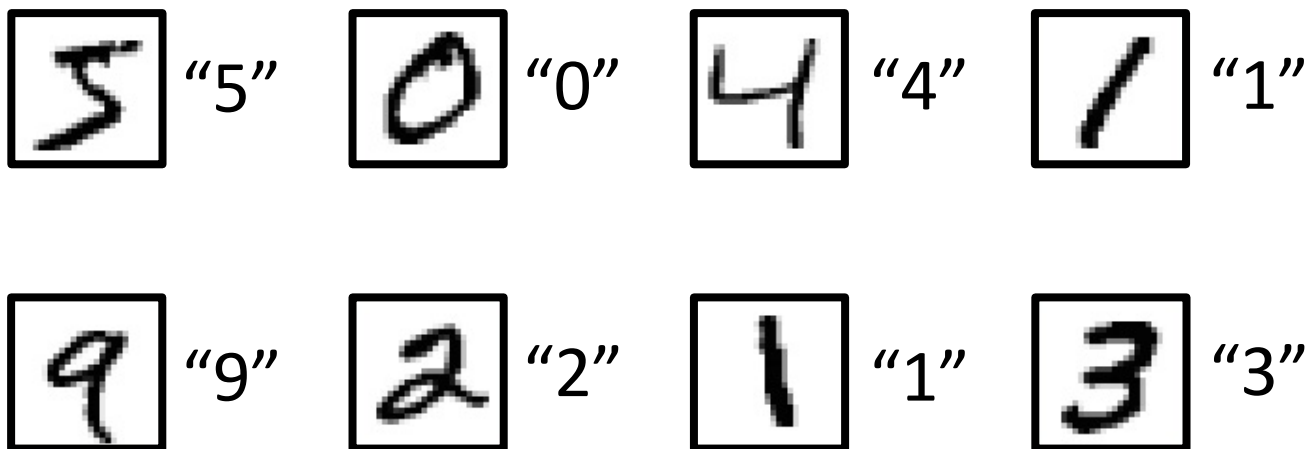
Input x_2 has the maximum value

How to let the neural network achieve this



Training Data

- Preparing training data: images and their labels



Using the training data to find
the network parameters.

Formalize

- Preparing training data: images and their labels

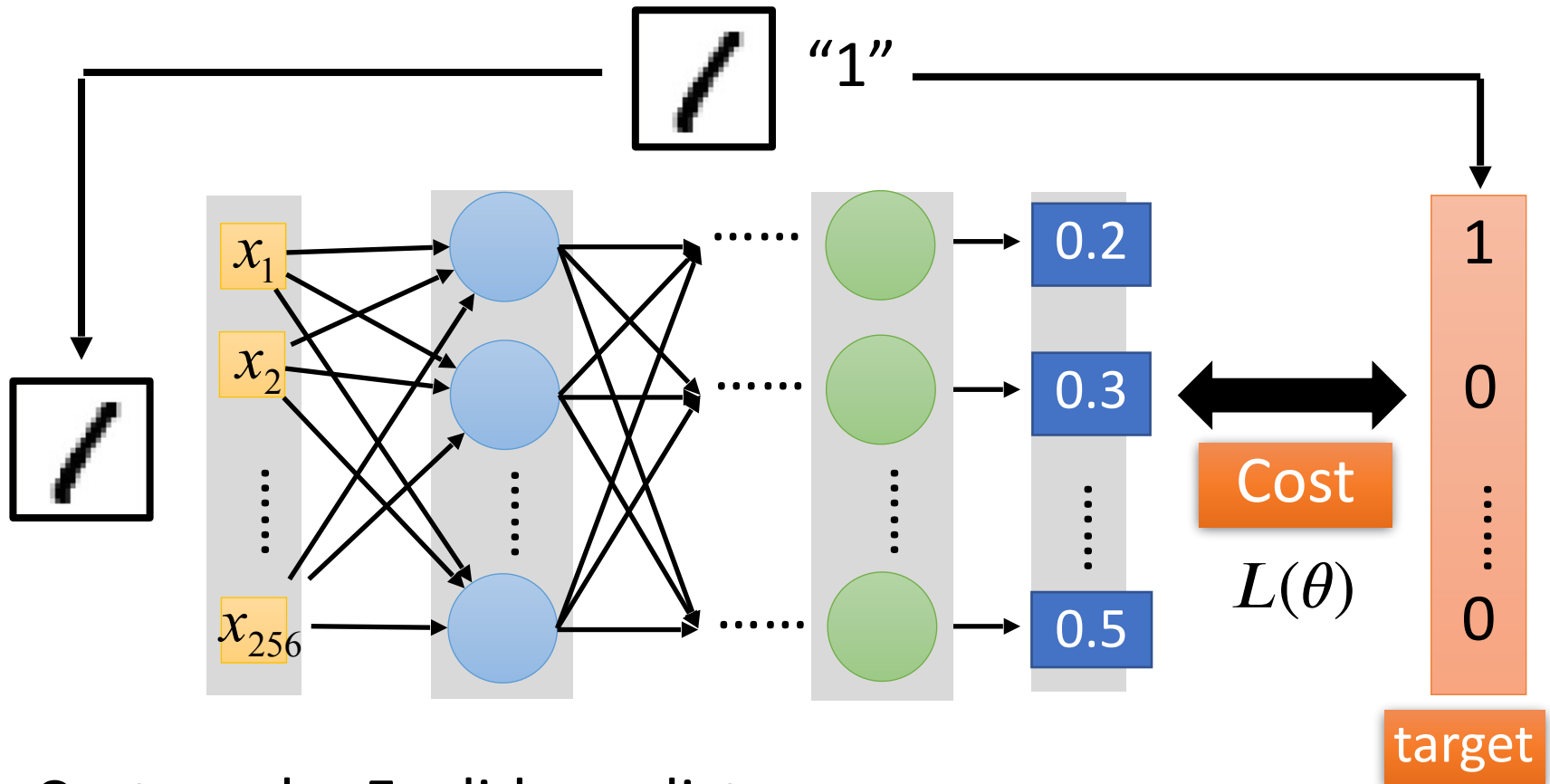
x :



$$y_0^{gt} = 0, \dots, y_4^{gt} = 0, y_5^{gt} = 1, y_6^{gt} = 0, \dots, y_9^{gt} = 0$$

Cost

Given a set of network parameters θ , each example has a cost value.



Cost can be Euclidean distance or cross entropy of the network output and target

Soft-Entropy Loss

The predicted score of groundtruth label category is larger than other categories:

$$y_{label} > y_j \text{ for any } j \neq label$$

How to set up a loss for this goal?

Soft-Entropy Loss

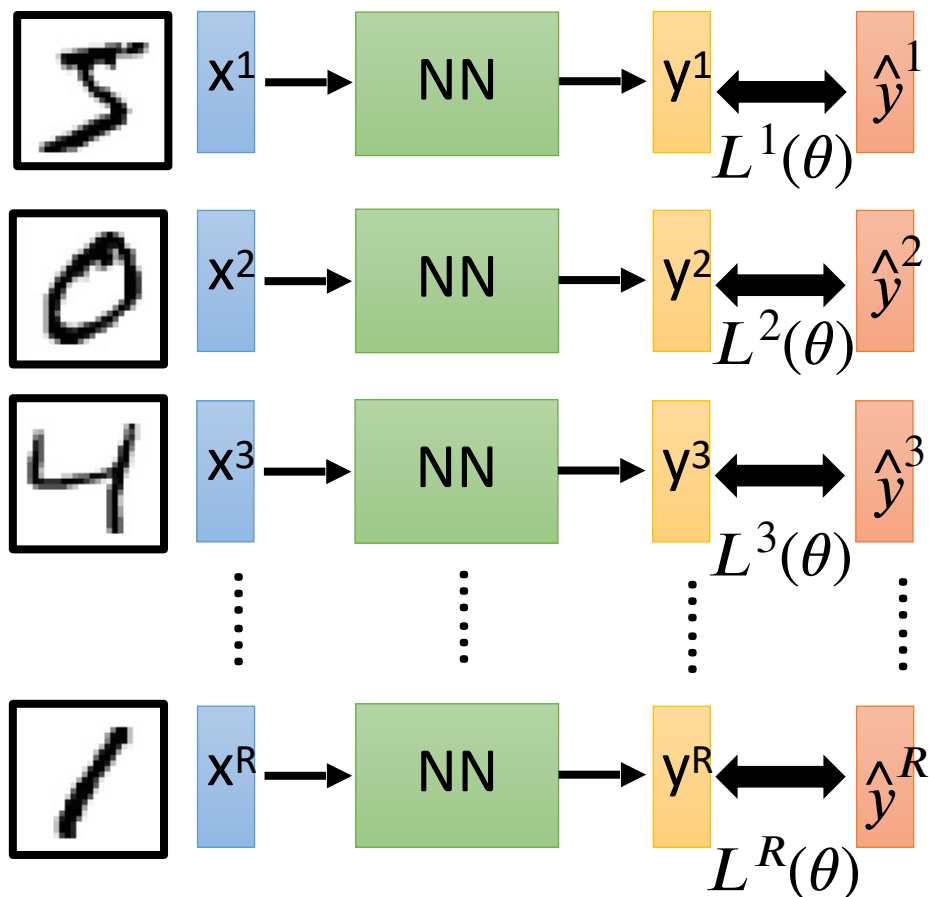
Let $y_{label}(x; \theta) = \frac{e^{f(x; \theta)_{label}}}{\sum_j e^{f(x; \theta)_j}}$

We minimize the loss

$$L(\theta) = -\log y_{label}(\theta)$$

Total Cost

For all training data ...



Total Cost:

$$C(\theta) = \sum_{r=1}^R L^r(\theta)$$

How bad the network parameters θ is on this task

Find the network parameters θ^* that minimize this value

Gradient Descent


Error Surface

Assume there are only two parameters w_1 and w_2 in a network.


$$\theta = \{w_1, w_2\}$$

Randomly pick a starting point θ^0

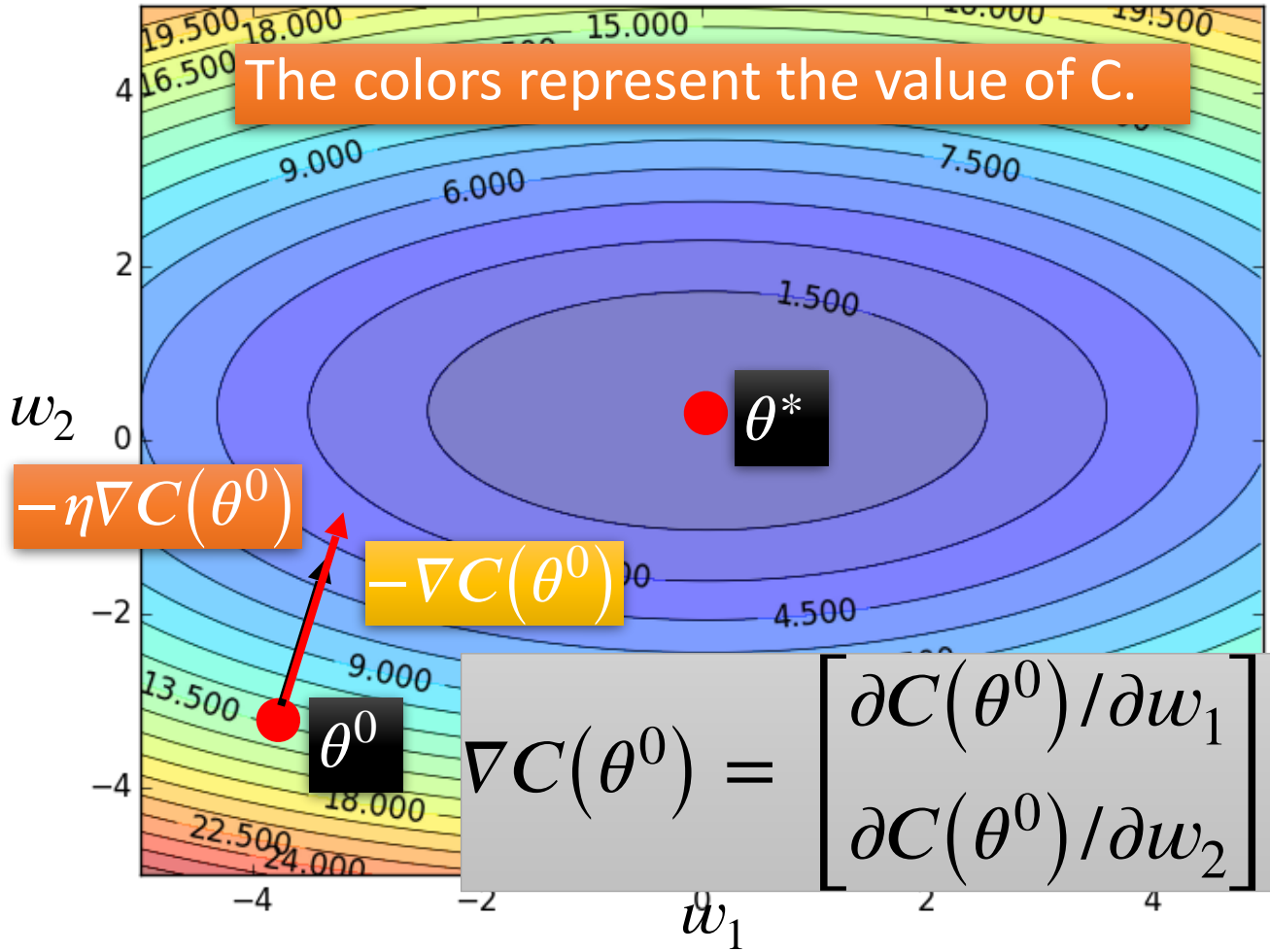
Compute the negative gradient at θ^0

 $-\nabla C(\theta^0)$

Times the learning rate η

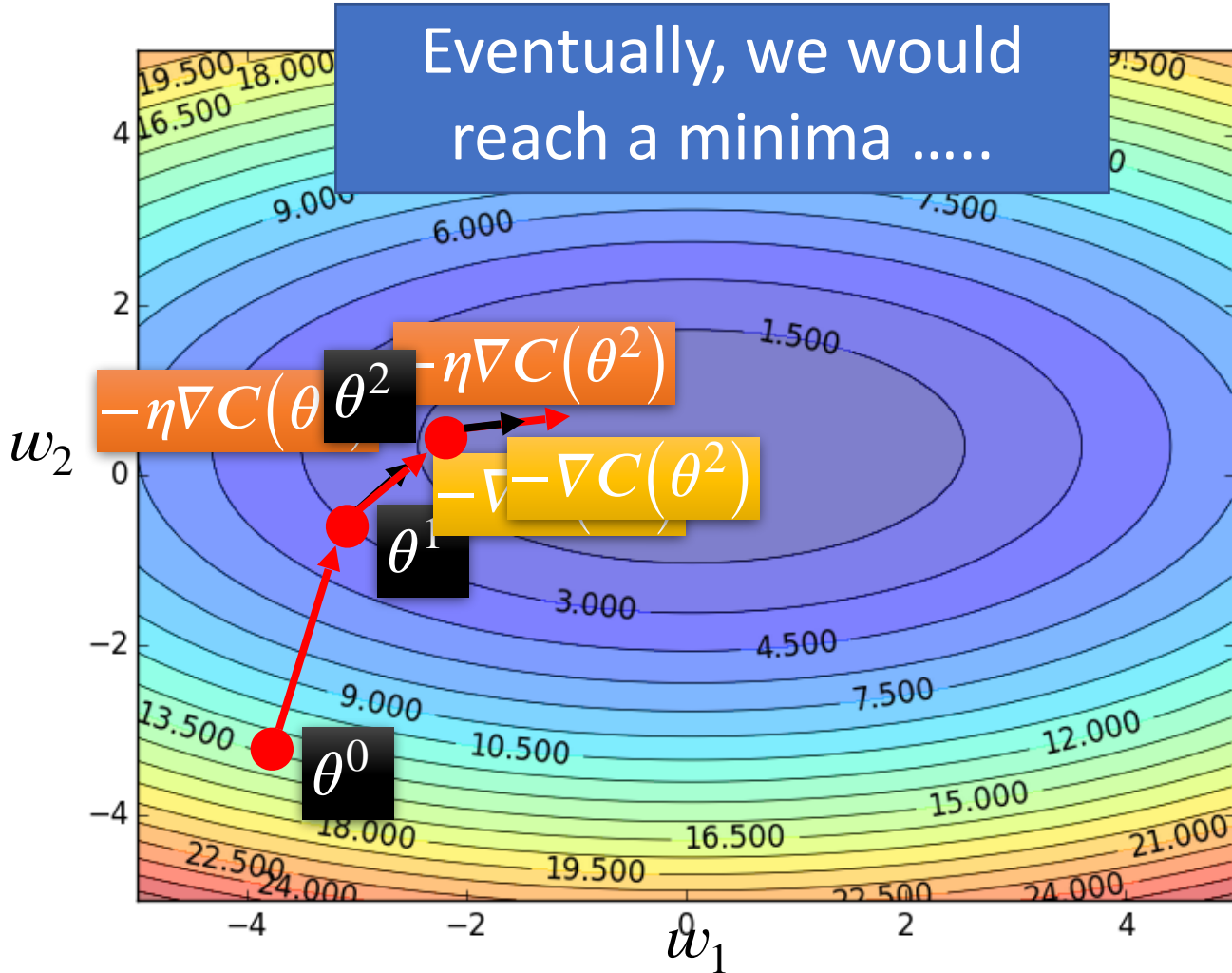
 $-\eta \nabla C(\theta^0)$

The colors represent the value of C.



$$\nabla C(\theta^0) = \begin{bmatrix} \partial C(\theta^0) / \partial w_1 \\ \partial C(\theta^0) / \partial w_2 \end{bmatrix}$$

Gradient Descent



Randomly pick a starting point θ^0

Compute the negative gradient at θ^0

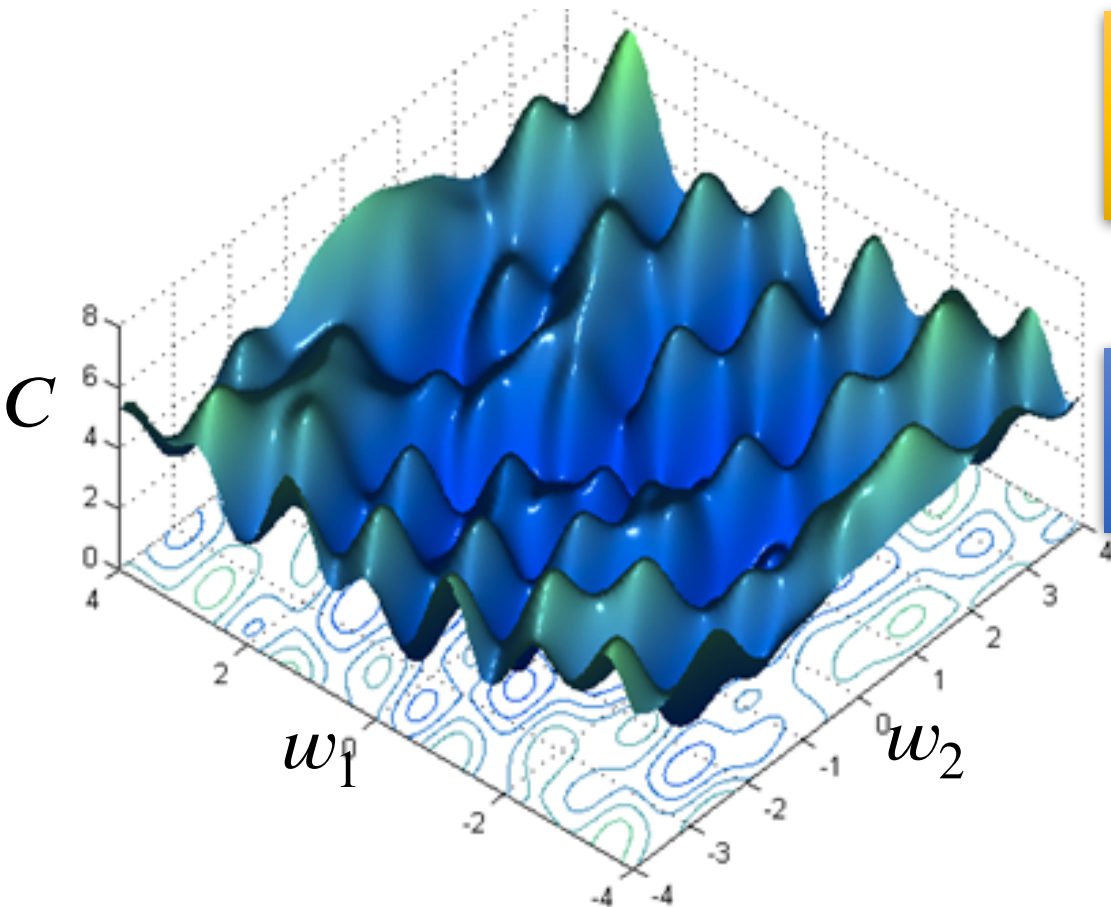
$\rightarrow -\nabla C(\theta^0)$

Times the learning rate η

$\rightarrow -\eta \nabla C(\theta^0)$

Local Minima

- Gradient descent never guarantee global minima



Different initial
point θ^0

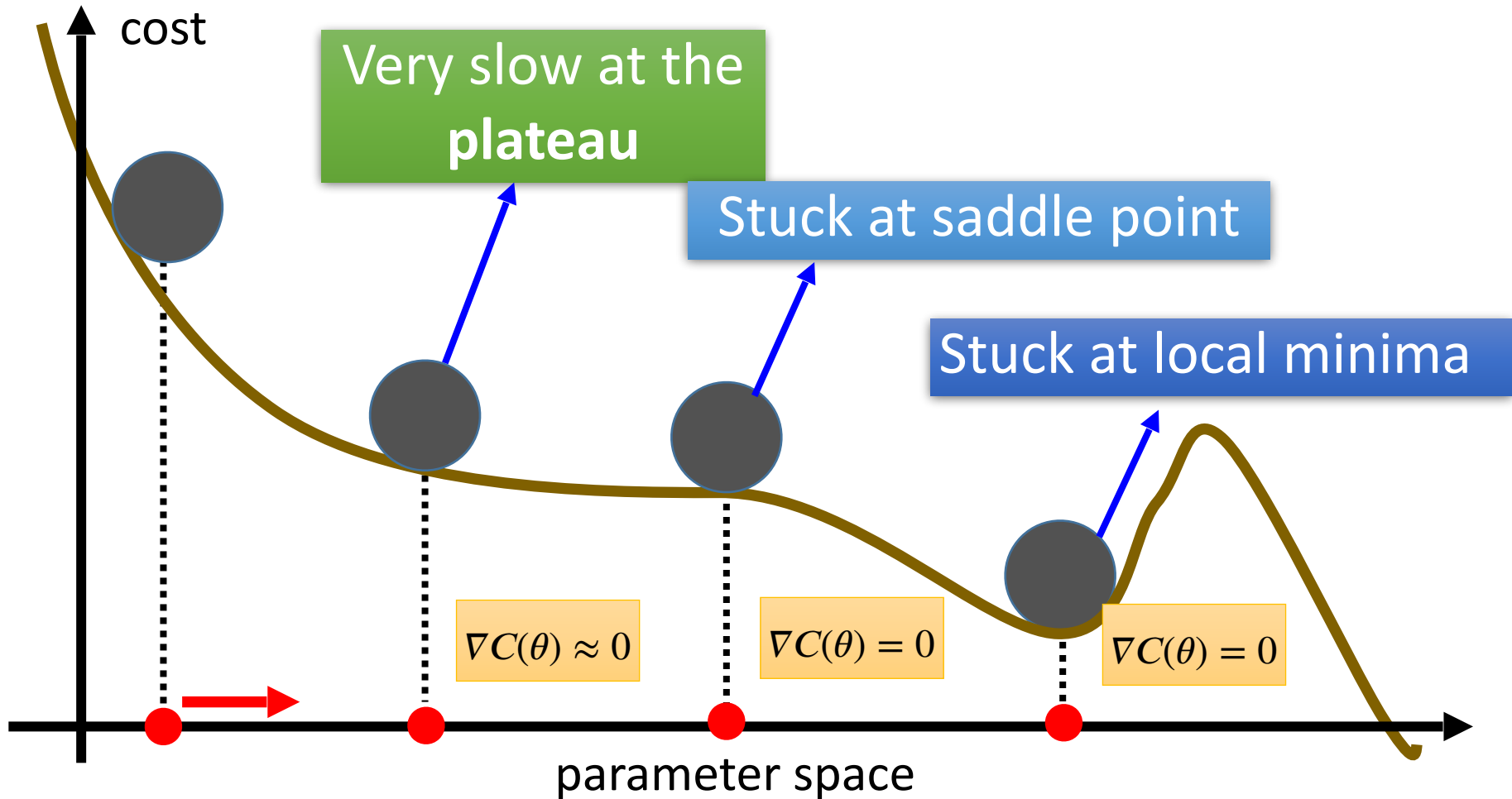


Reach different minima,
so different results

Who is Afraid of Non-Convex
Loss Functions?

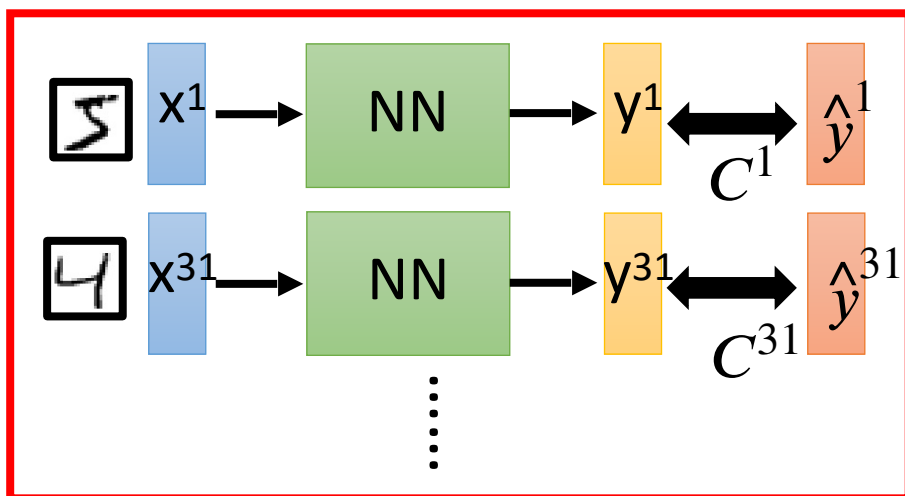
[http://videolectures.net/
eml07_lecun_wia/](http://videolectures.net/eml07_lecun_wia/)

Besides local minima

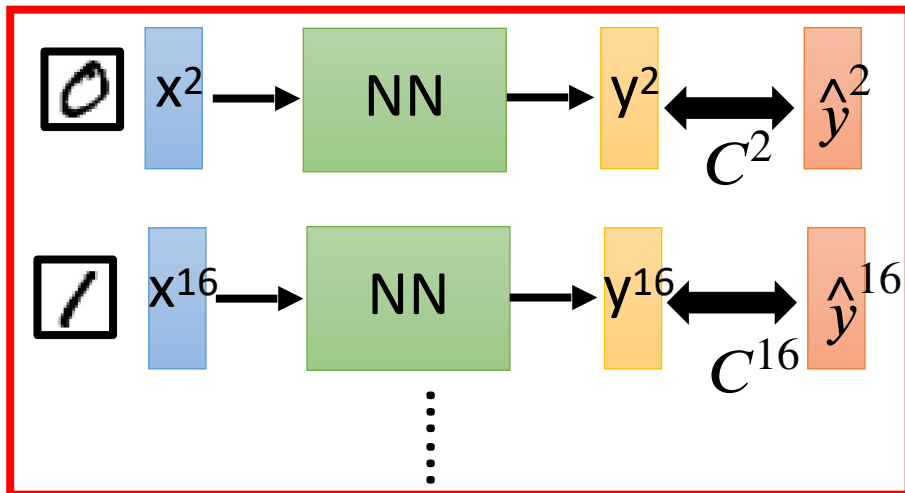


Stochastic Gradient Descent (SGD)

Mini-batch



Mini-batch



➤ Randomly initialize θ^0

➤ Pick the 1st batch

$$C = C^1 + C^{31} + \dots$$

$$\theta^1 \leftarrow \theta^0 - \eta \nabla C(\theta^0)$$

➤ Pick the 2nd batch

$$C = C^2 + C^{16} + \dots$$

$$\theta^2 \leftarrow \theta^1 - \eta \nabla C(\theta^1)$$

⋮

➤ Until all mini-batches have been picked

one epoch

Repeat the above process

Backpropagation

- A network can have millions of parameters.
 - Backpropagation is the way to compute the gradients efficiently (not today)
 - Ref: http://speech.ee.ntu.edu.tw/~tlkagk/courses/MLDS_2015_2/Lecture/DNN%20backprop.ecm.mp4/index.html
- Many toolkits can compute the gradients automatically

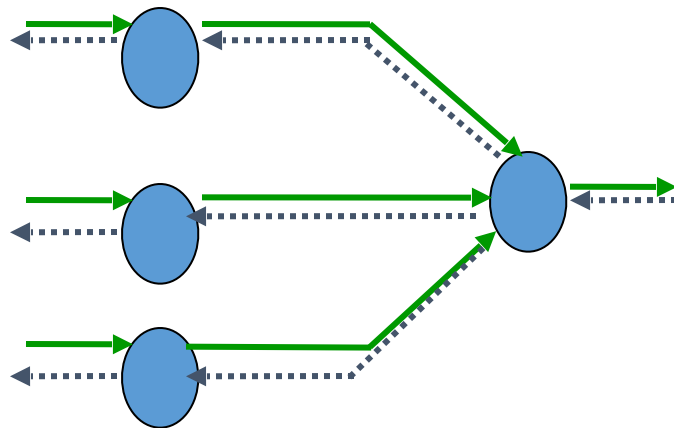
theano



Ref: http://speech.ee.ntu.edu.tw/~tlkagk/courses/MLDS_2015_2/Lecture/Theano%20DNN.ecm.mp4/index.html

Back Propagation

- Back-propagation training algorithm



*Network activation
Forward Step*

*Error propagation
Backward Step*

- Backprop adjusts the weights of the NN in order to minimize the network total error.

Why Deep?

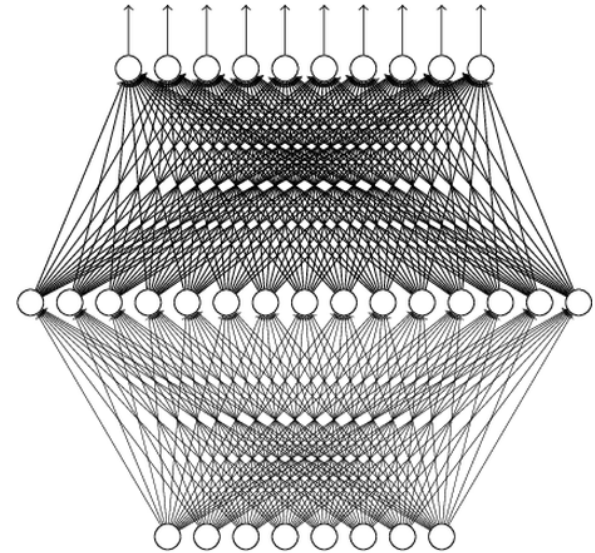
Universality Theorem

Any continuous function f

$$f : \mathbb{R}^N \rightarrow \mathbb{R}^M$$

Can be realized by a network
with one hidden layer

(given **enough** hidden
neurons)



Reference for the reason:

[http://](http://neuralnetworksanddeeplearning.com/chap4.html)

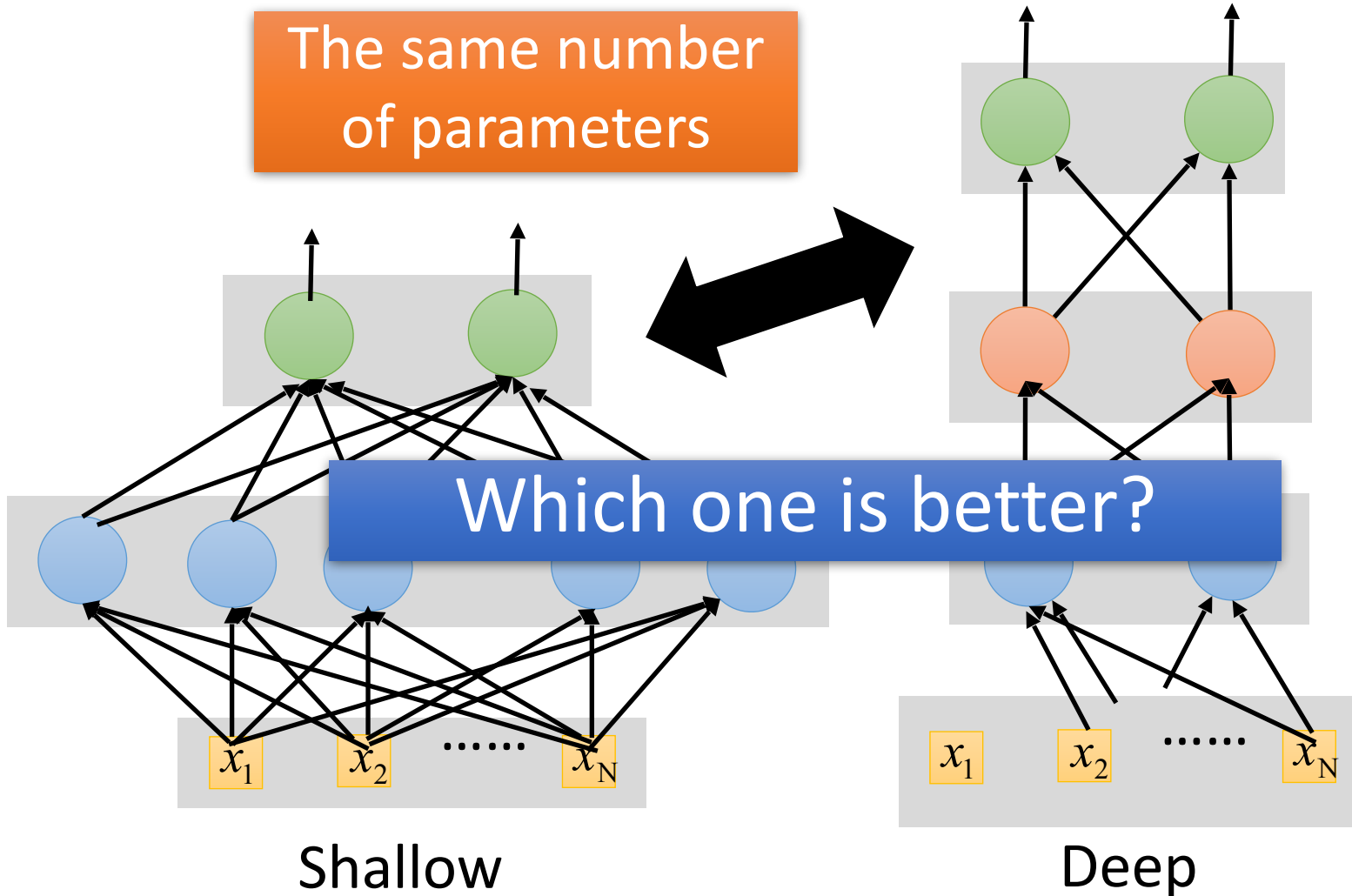
neuralnetworksanddeeplearning.com/chap4.html

The Unreasonable Effectiveness of Gradient Descent

- While the loss function for neural networks is highly non-convex, empirically (and theoretically), we can show that, with many hidden neurons, the value of local minima are almost as **small** as the global minimum

Then why “Deep” neural network not “Fat” neural network?

Fat + Short v.s. Thin + Tall



Fat + Short v.s. Thin + Tall

“Why deep” is a very “deep” question!

No simple answer yet, even no fully
convincing answer yet!

Statistical View of Machine Learning

We start from understanding some simple classifiers,
to draw inspiration for understanding neural networks!

Example Task: Image Classification



This image by Nikita is licensed under [CC-BY 2.0](https://creativecommons.org/licenses/by/2.0/)

(assume given set of discrete labels)
{dog, cat, truck, plane, ...}

→ cat

An image classifier

```
def classify_image(image):  
    # Some magic here?  
    return class_label
```

Unlike e.g. sorting a list of numbers,

no obvious way to hard-code the algorithm for recognizing a cat, or other classes.

Machine Learning: Data-Driven Approach

1. Collect a dataset of images and labels
2. Use Machine Learning to train a classifier
3. Evaluate the classifier on new images

Example training set

```
def train(images, labels):  
    # Machine learning!  
    return model
```

```
def predict(model, test_images):  
    # Use model to predict labels  
    return test_labels
```

airplane



automobile



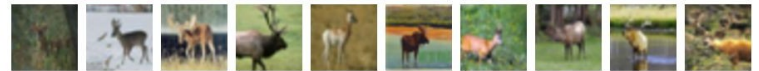
bird



cat



deer



First classifier: Nearest Neighbor

```
def train(images, labels):  
    # Machine learning!  
    return model
```



Memorize all
data and labels

```
def predict(model, test_images):  
    # Use model to predict labels  
    return test_labels
```



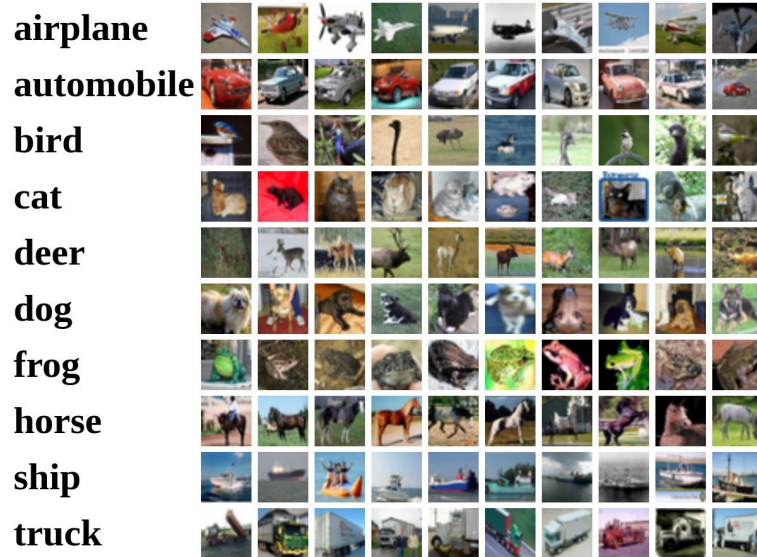
Predict the label
of the most similar
training image

Example Dataset: CIFAR10

10 classes

50,000 training images

10,000 testing images

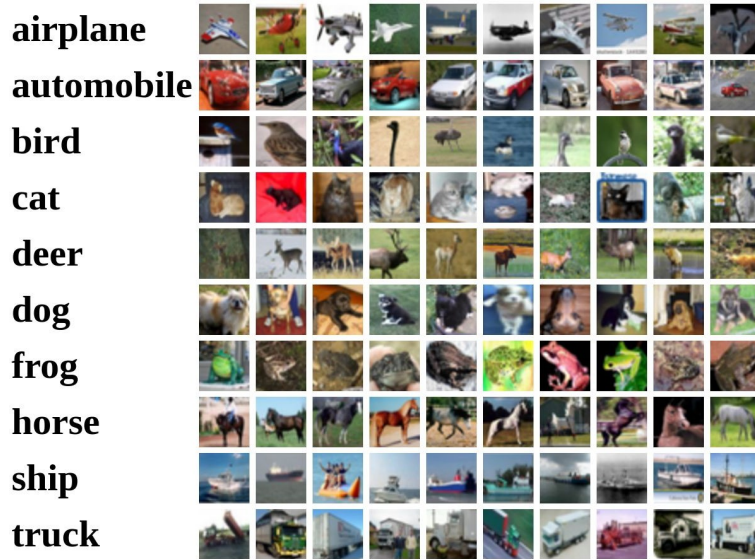


Example Dataset: CIFAR10

10 classes

50,000 training images

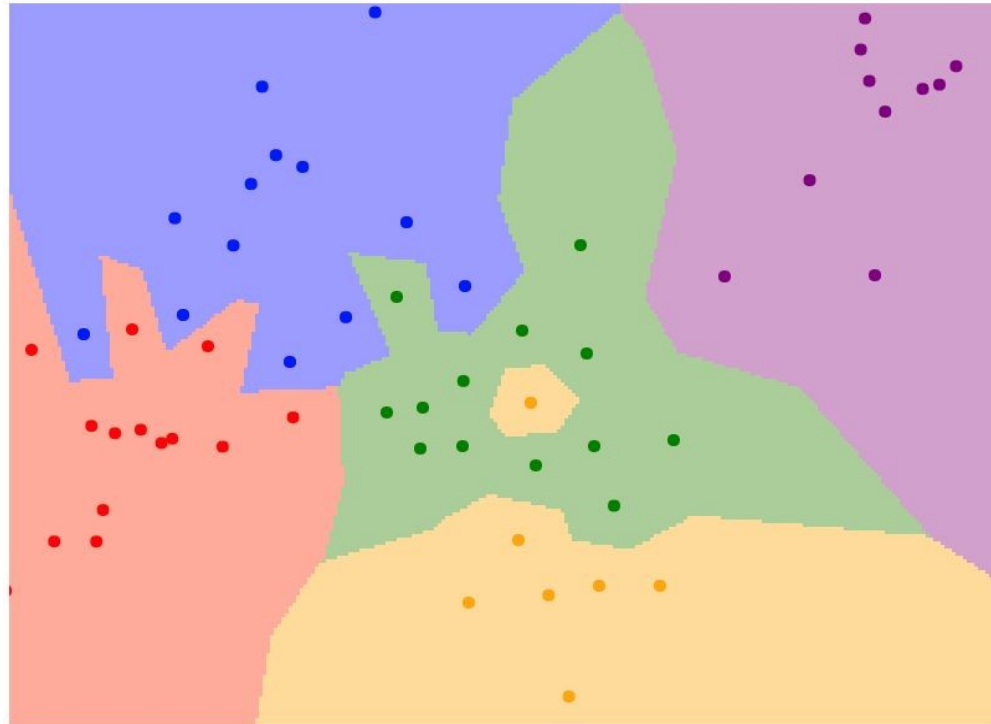
10,000 testing images



Test images and nearest neighbors

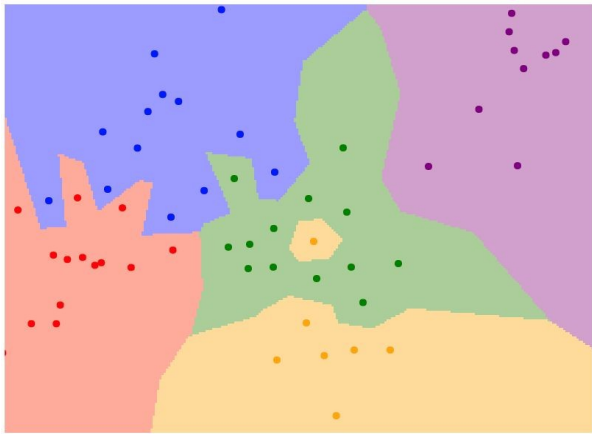


What does Nearest Neighbor look like?

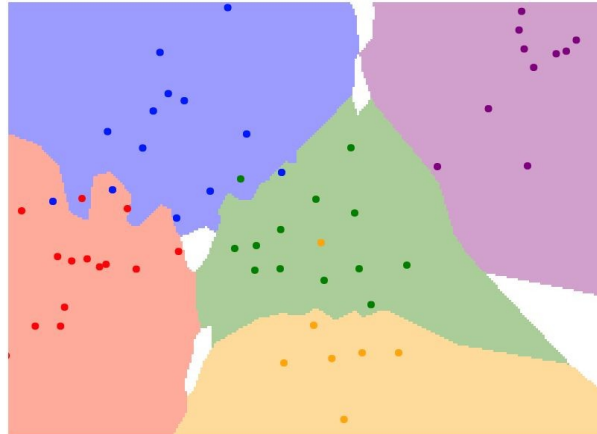


K-Nearest Neighbors

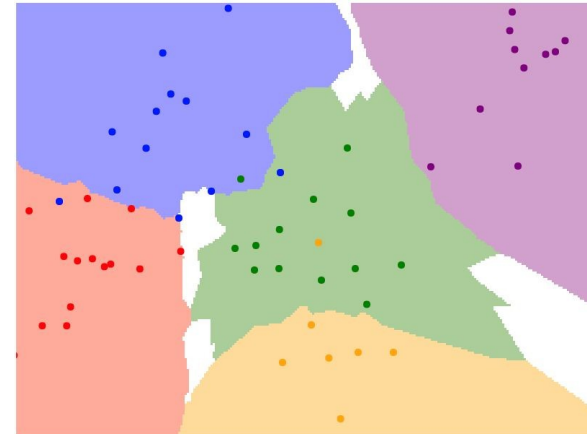
Instead of copying label from nearest neighbor,
take **majority vote** from K closest points



$K = 1$



$K = 3$

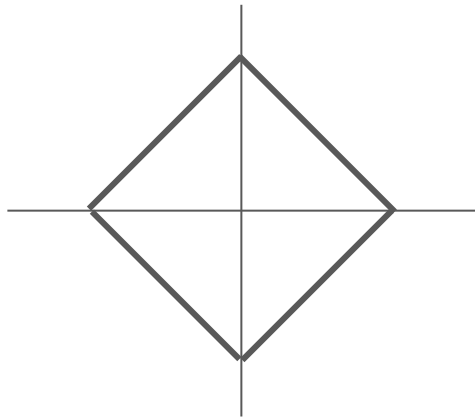


$K = 5$

K-Nearest Neighbors: Distance Metric

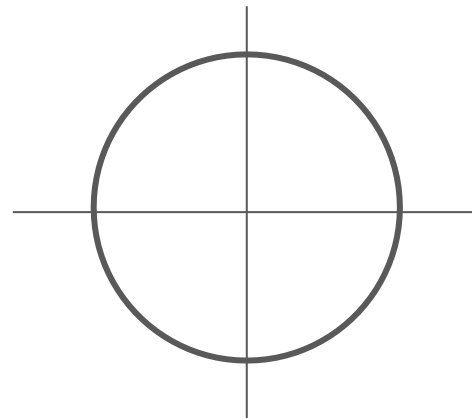
L1 (Manhattan) distance

$$d_1(I_1, I_2) = \sum_p |I_1^p - I_2^p|$$



L2 (Euclidean) distance

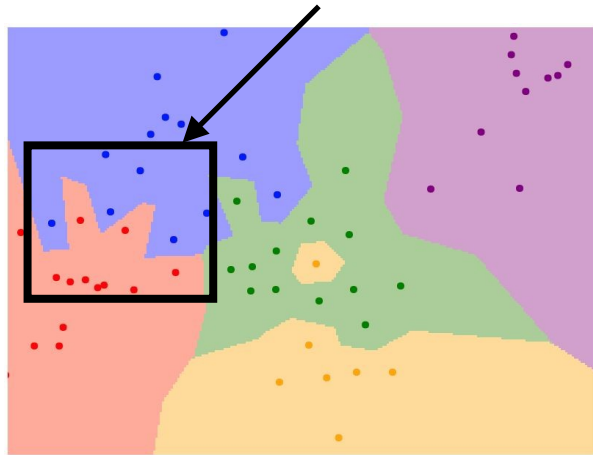
$$d_2(I_1, I_2) = \sqrt{\sum_p (I_1^p - I_2^p)^2}$$



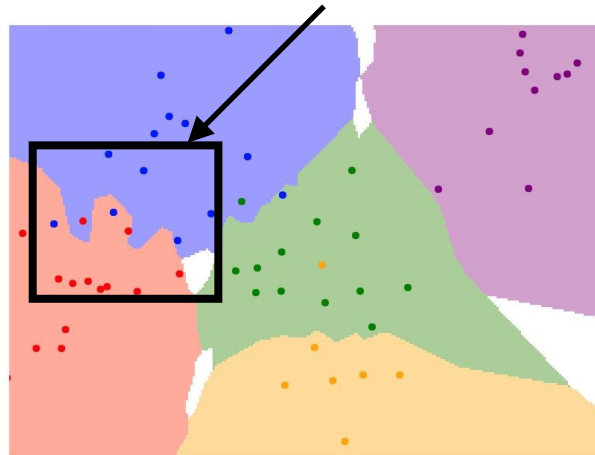
Hyperparameters

- What is the best value of k to use? What is the best distance to use?
- These are hyperparameters: choices about the algorithm that we set rather than learn
- The deep v.s. fat choice for neural networks is similarly a choice of “algorithms”

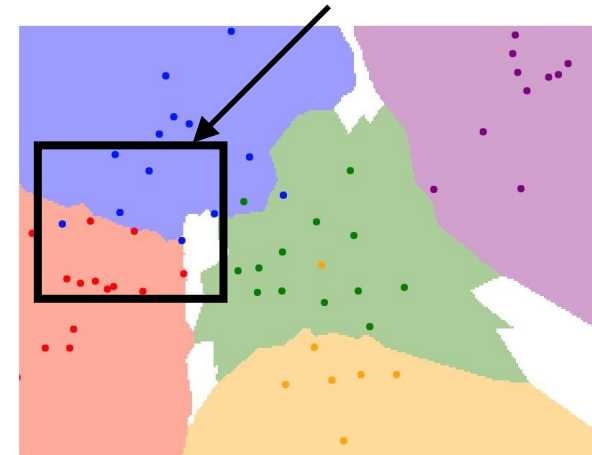
K-Nearest Neighbors



$K = 1$



$K = 3$



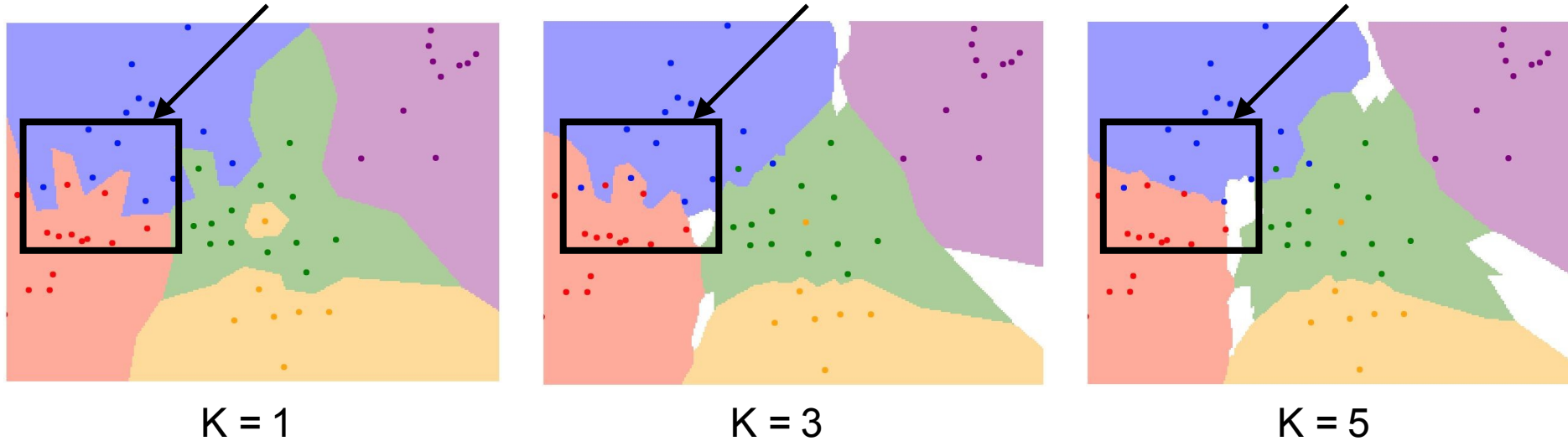
$K = 5$

Observations:

Small K (e.g., $K=1$): every sample matters, sophisticated boundary

Large K (e.g., $K=5$): voting finds the consensus in the neighborhood, simpler boundary

K-Nearest Neighbors



Observations:

Small K (e.g., $K=1$): every sample matters, sophisticated boundary, **high model complexity**

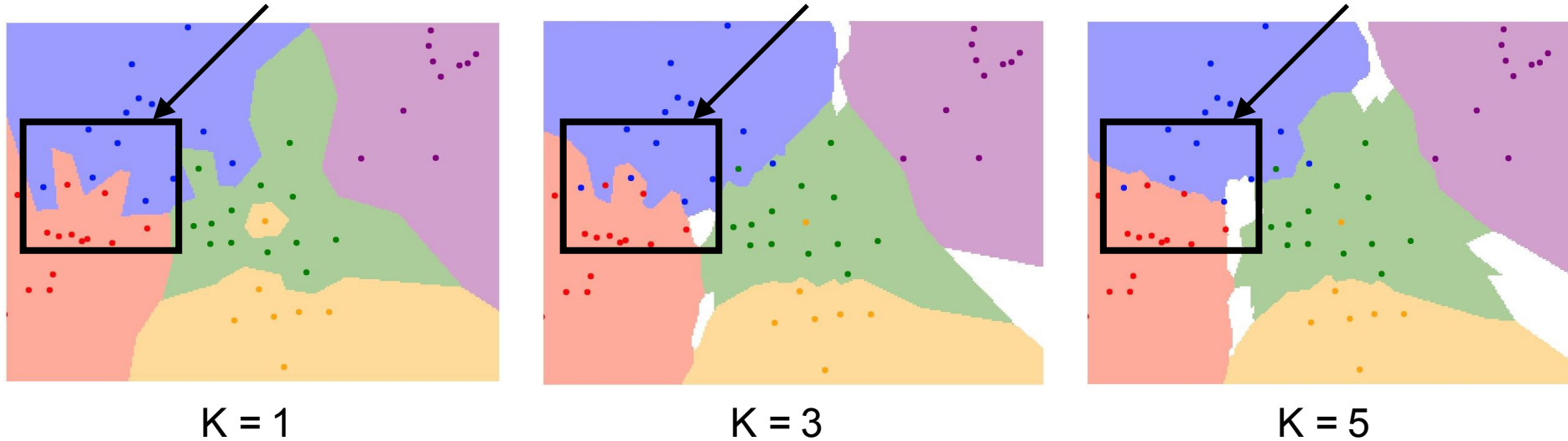
Large K (e.g., $K=5$): voting finds the consensus in the neighborhood, simpler boundary, **low model complexity**

Bias and Variance

- Bias – error caused because the model lacks the ability to represent the (complex) concept
- Variance – error caused because the learning algorithm overreacts to small changes (noise) in the training data

$$\text{TotalLoss} = \text{Bias} + \text{Variance} (+ \text{noise})$$

K-Nearest Neighbors



Which one has higher bias? higher variance?

- Bias – error caused because the model lacks the ability to represent the (complex) concept
- Variance – error caused because the learning algorithm overreacts to small changes (noise) in the training data

The Power of a Model Building Process

Weaker Modeling Process (higher bias)

- Simple Model (e.g. linear, large K in KNN)
- Small Feature Set (e.g. few neurons)
- Constrained Search (e.g. few iterations of gradient descent)

More Powerful Modeling Process (higher variance)

- Complex Model (e.g. networks, small K in KNN)
- Large Feature Set (e.g. many neurons)
- Unconstrained Search (e.g. exhaustive search)

Overfitting v.s. Underfitting

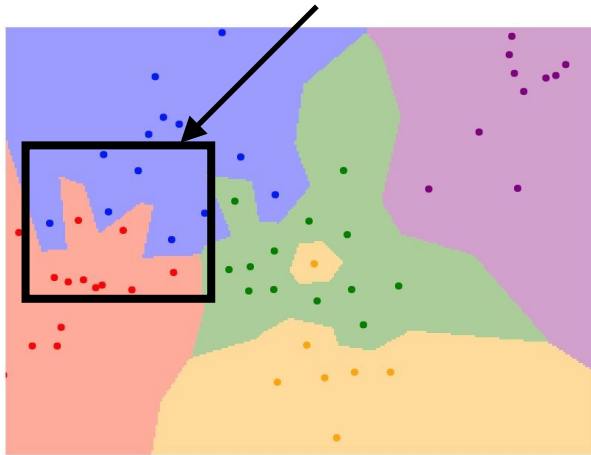
Overfitting

- Fitting the data too well
 - Features are noisy / uncorrelated to concept
 - Modeling process very sensitive (powerful)
 - Too much search

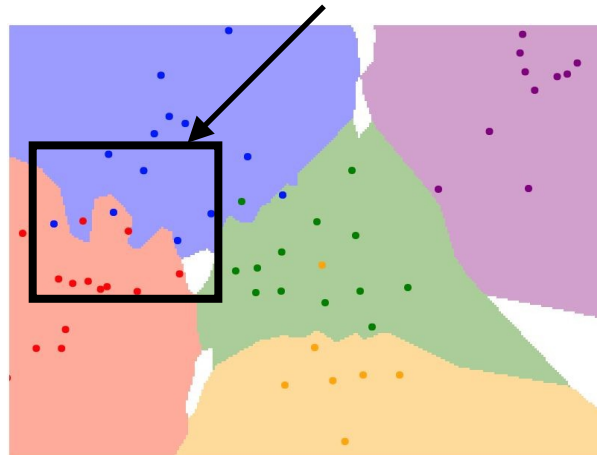
Underfitting

- Learning too little of the true concept
 - Features don't capture concept
 - Too much bias in model
 - Too little search to fit model

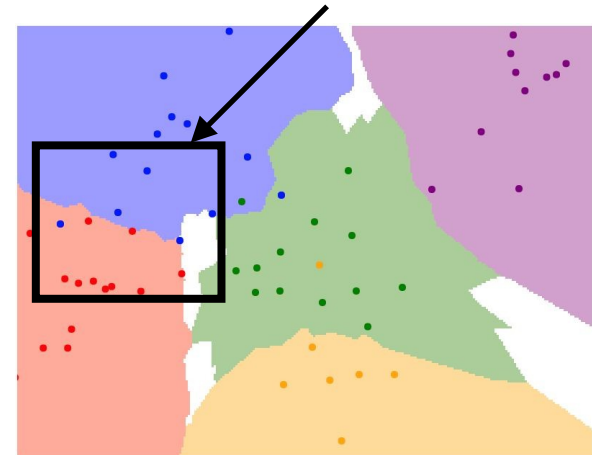
K-Nearest Neighbors



$K = 1$



$K = 3$



$K = 5$

Which one tends to overfit? to underfit?

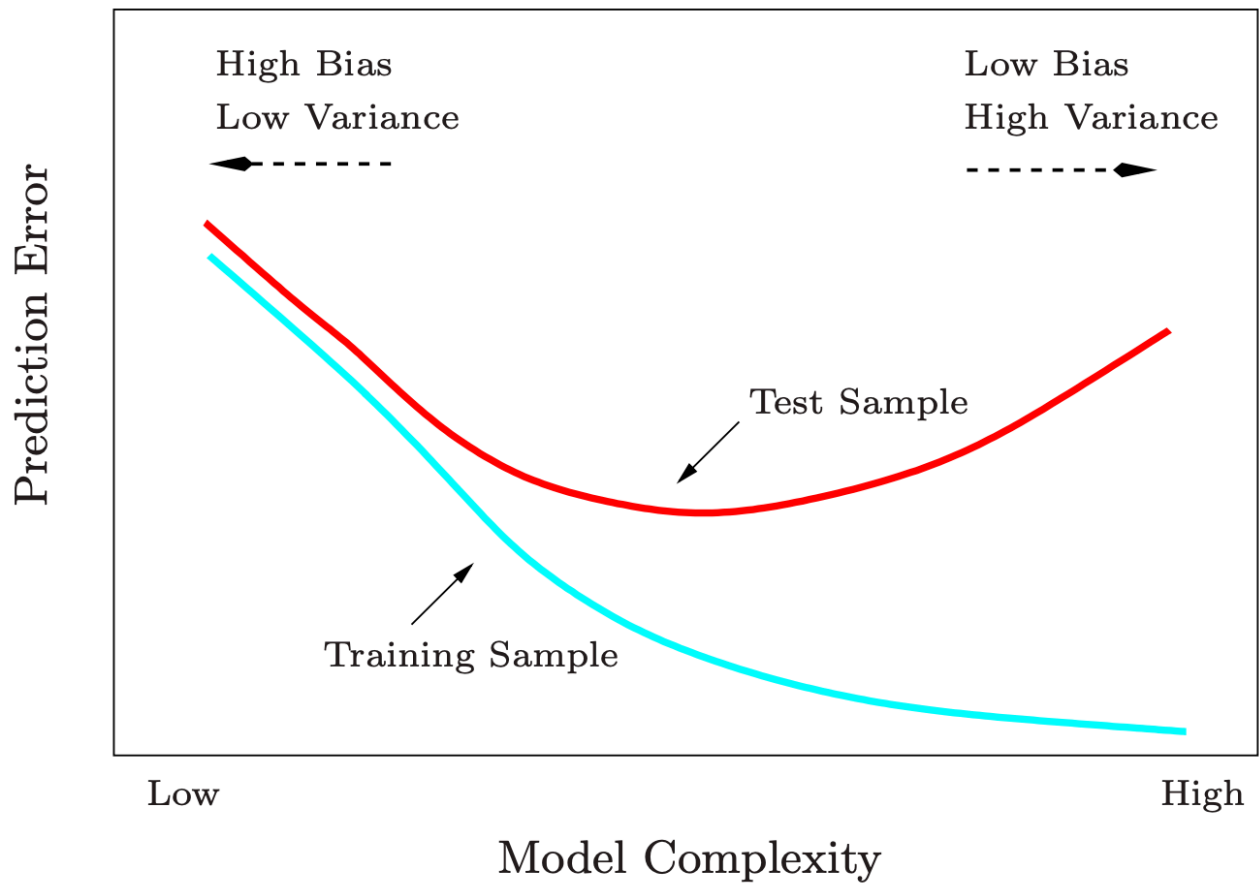


FIGURE 2.11. *Test and training error as a function of model complexity.*

Summary of Overfitting and Underfitting

- Bias / Variance tradeoff a primary challenge in machine learning
- Internalize: More powerful modeling is not always better
- Learn to identify overfitting and underfitting
- Tuning parameters & interpreting output correctly is key

Back to Neural Networks

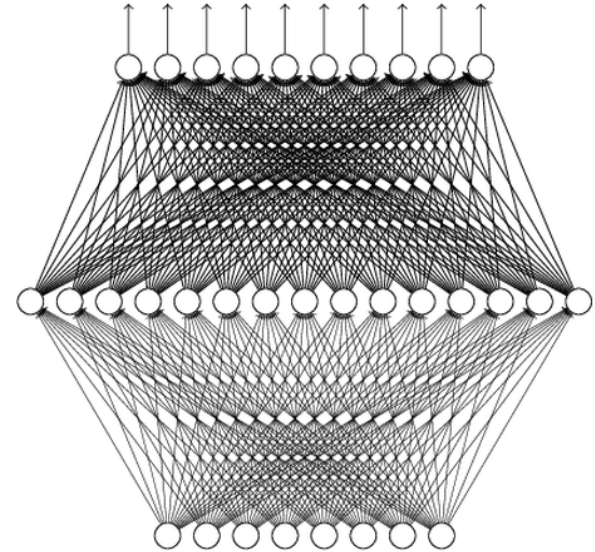
Recap: Universality Theorem

Any continuous function f

$$f : \mathbb{R}^N \rightarrow \mathbb{R}^M$$

Can be realized by a network
with one hidden layer

(given **enough** hidden
neurons)



Reference for the reason:

[http://](http://neuralnetworksanddeeplearning.com/chap4.html)

neuralnetworksanddeeplearning.com/chap4.html

Universality is Not Enough

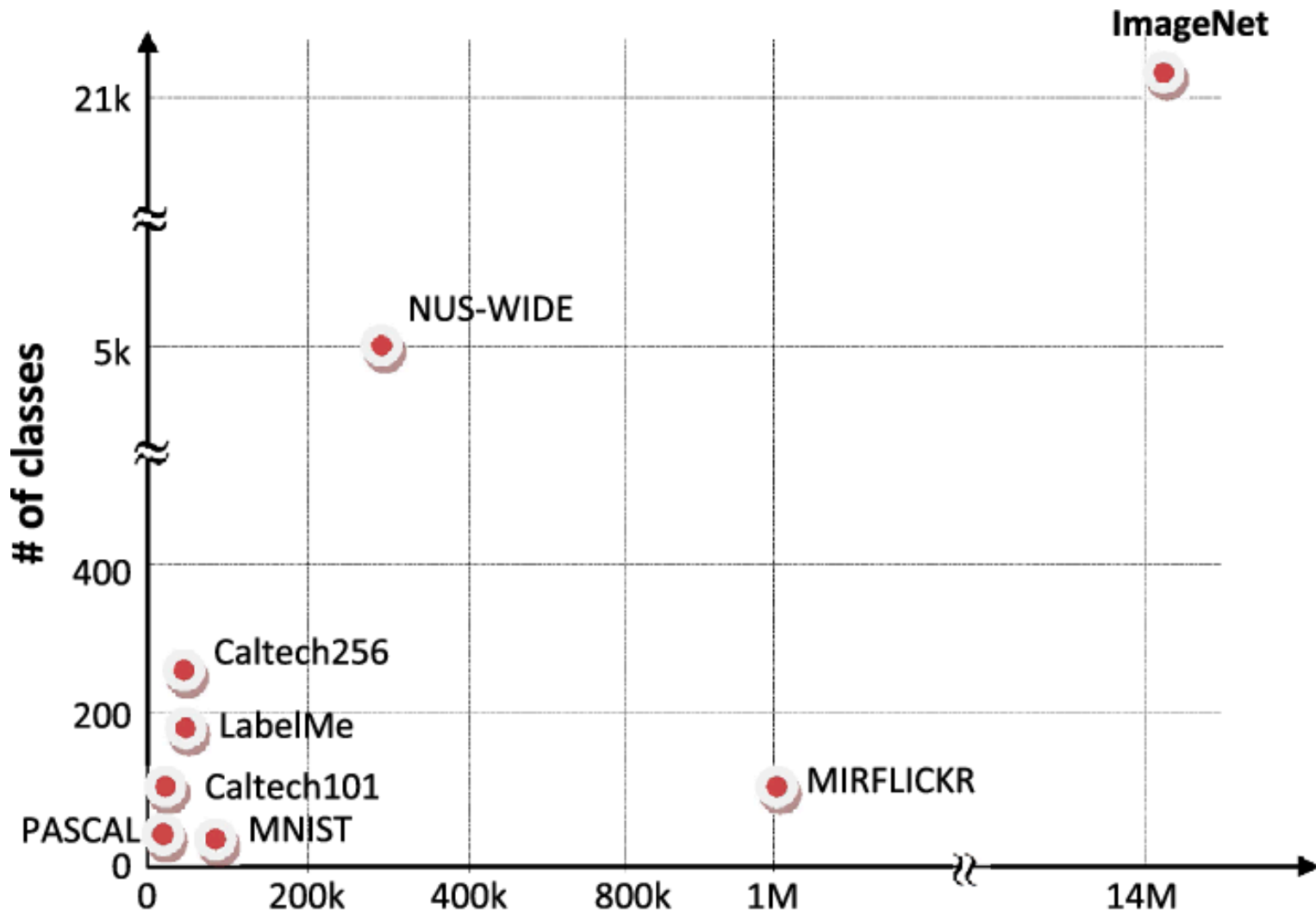
- Neural network has very high capacity (millions of parameters)
- By our basic knowledge of bias-variance tradeoff, so many parameters should imply very low bias, and very high variance. The test loss may not be small.
- Many efforts of deep learning are about mitigating overfitting!

Address Overfitting for NN

- Use larger training data set
- Design better network architecture

Address Overfitting for NN

- **Use larger training data set**
- Design better network architecture

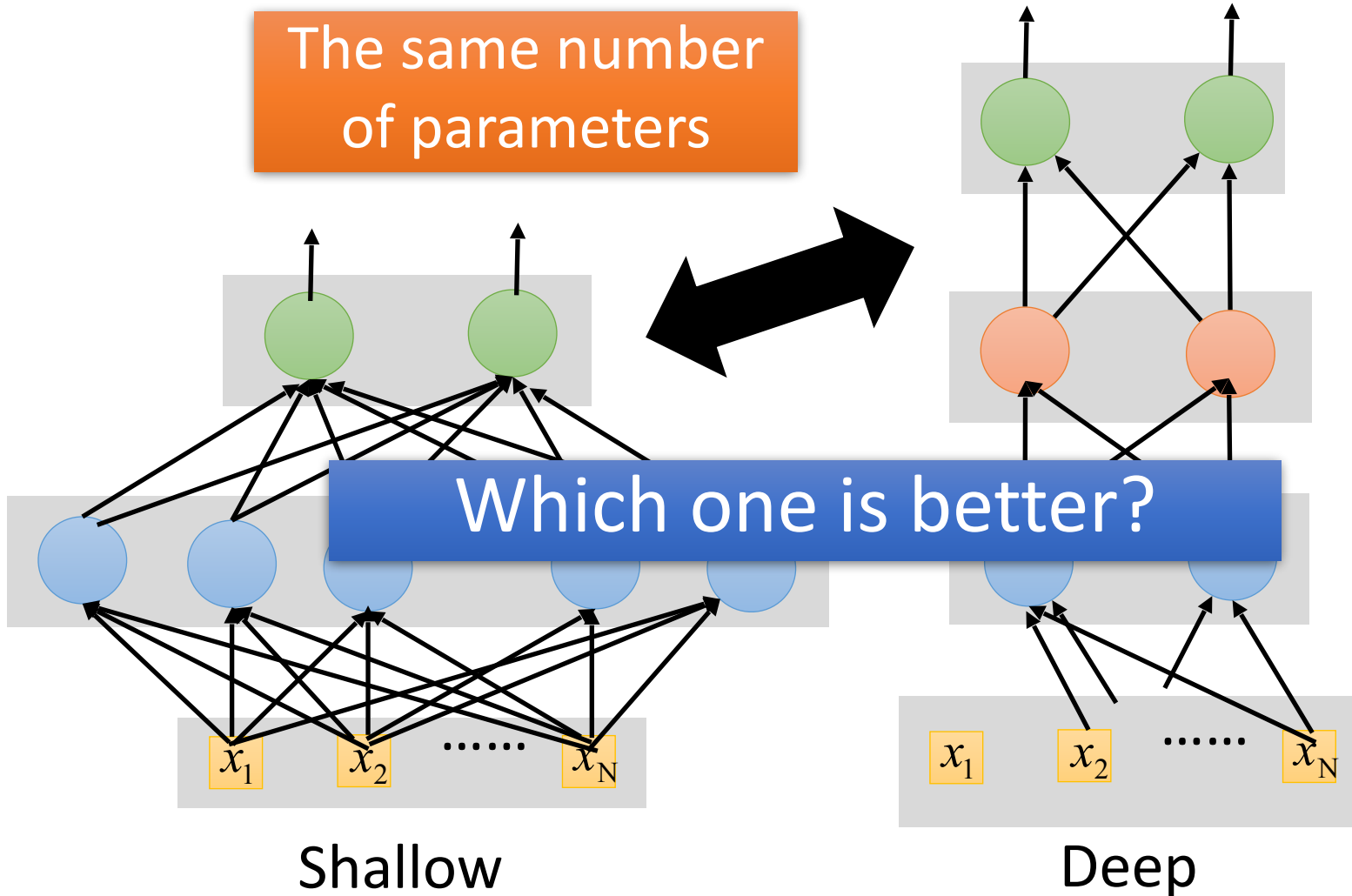


ImageNet Large Scale Visual Recognition Challenge
 Russakovsky, Deng, Su, et al. IJCV 2015

Address Overfitting for NN

- **Design better network architecture**
- Use larger training data set

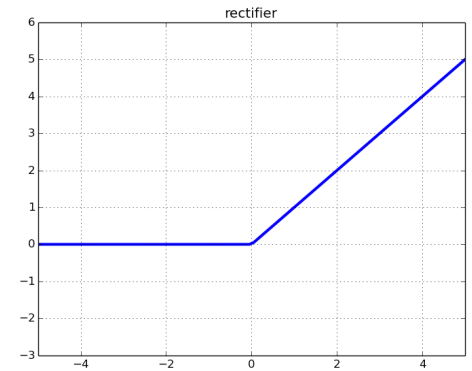
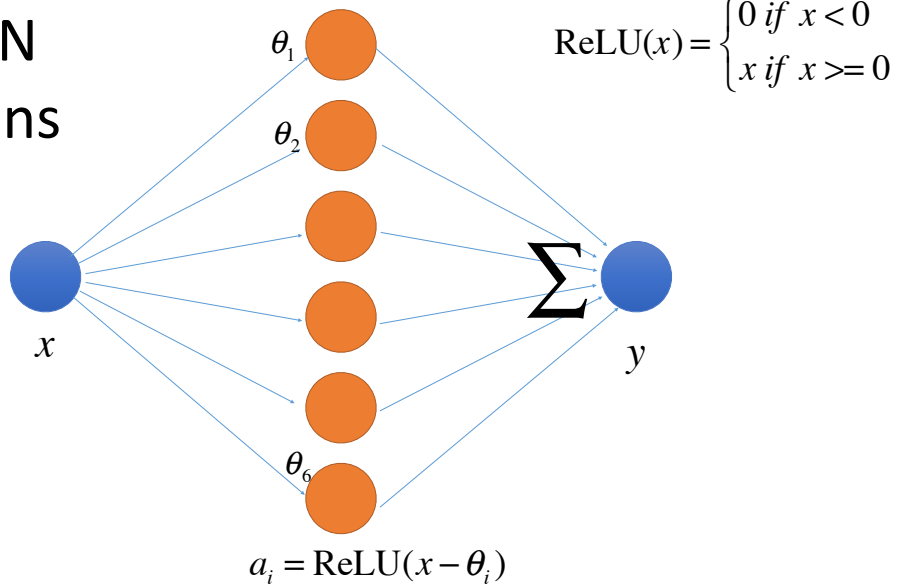
Fat + Short v.s. Thin + Tall



The Intuition behind Deep

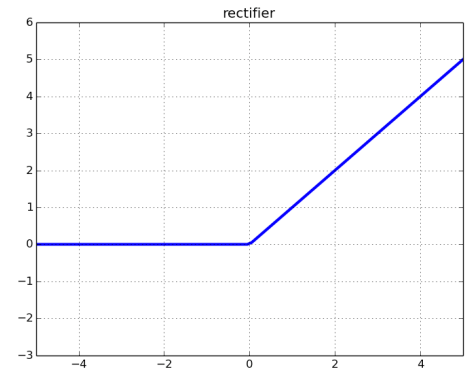
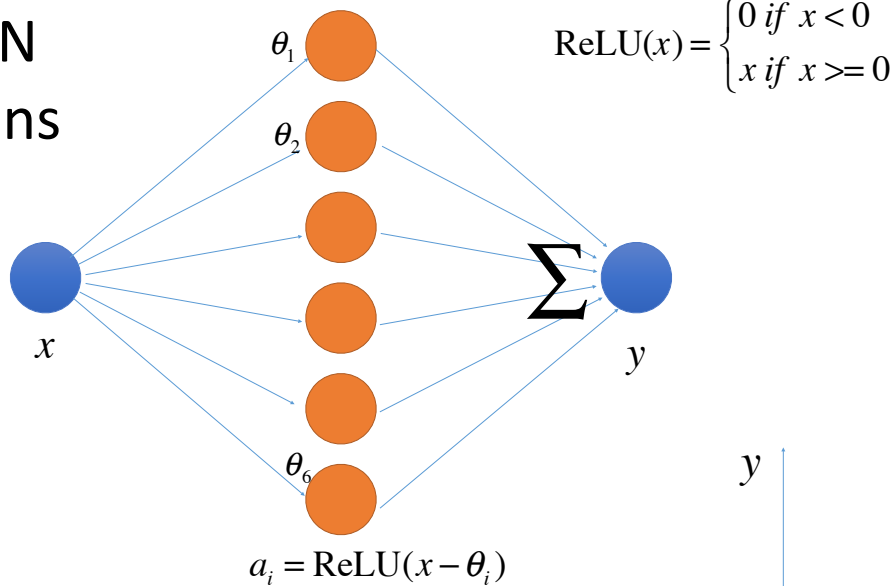
- To achieve the same representation power, we can use fewer neurons with a deeper architecture
- Fewer neurons risk less for overfitting (lacking rigor for this argument)

Fat + Short NN With 6 neurons

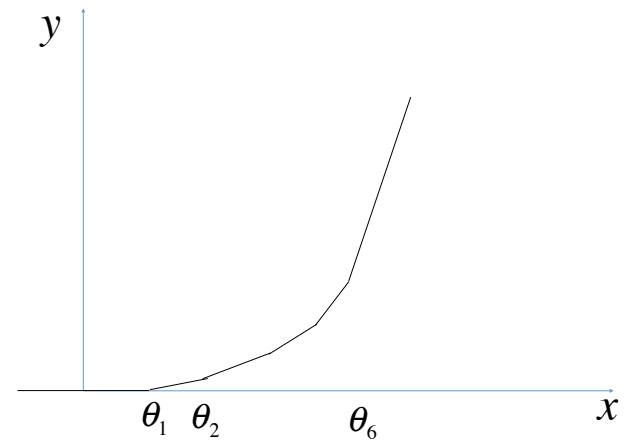


Assume $w_i = 1$ for all neurons, just learn the bias term b_i

Fat + Short NN With 6 neurons



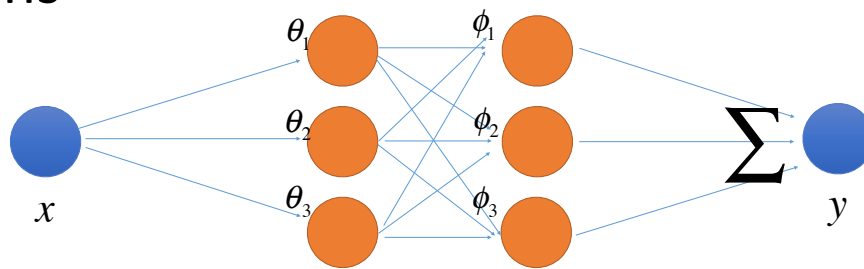
$$y = \sum_i \text{ReLU}(x - \theta_i)$$



piece-wise linear, 6 knots

Assume $w_i = 1$ for all neurons, just learn the bias term b_i

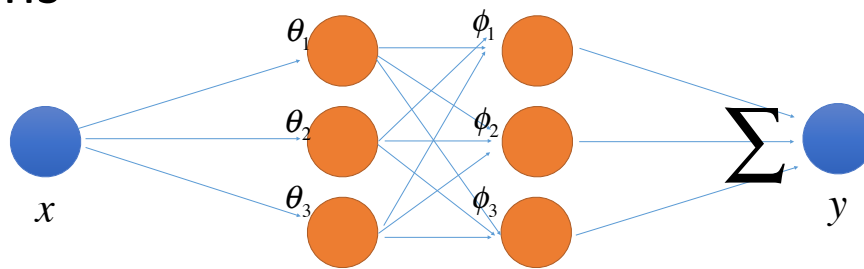
Thin + Tall NN With 6 neurons



$$a_{1,i} = \text{ReLU}(x - \theta_i) \quad a_{2,i} = \text{ReLU}(x - \phi_i) \quad y$$

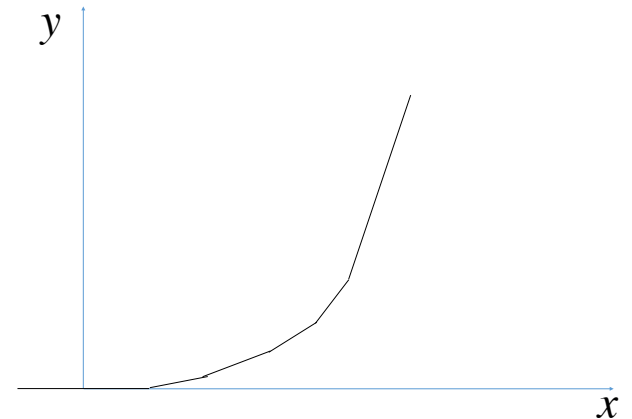
Assume $w_i = 1$ for all neurons, just learn the bias term b_i

Thin + Tall NN With 6 neurons



$$a_{1,i} = \text{ReLU}(x - \theta_i) \quad a_{2,i} = \text{ReLU}(x - \phi_i) \quad y$$

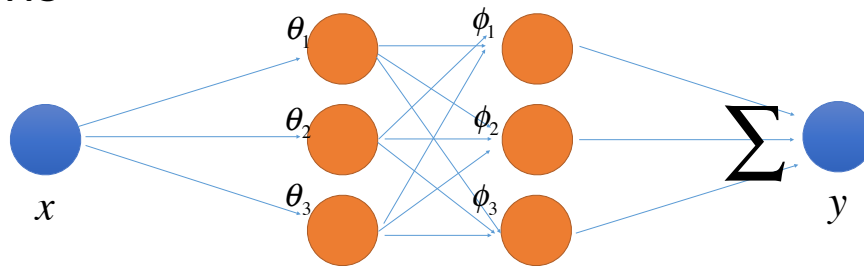
$$y = \sum_{1 \leq i \leq 3} \text{ReLU}([\sum_{1 \leq j \leq 3} \text{ReLU}(x - \theta_j)] - \phi_i)$$



piece-wise linear, can have **9 knots!**

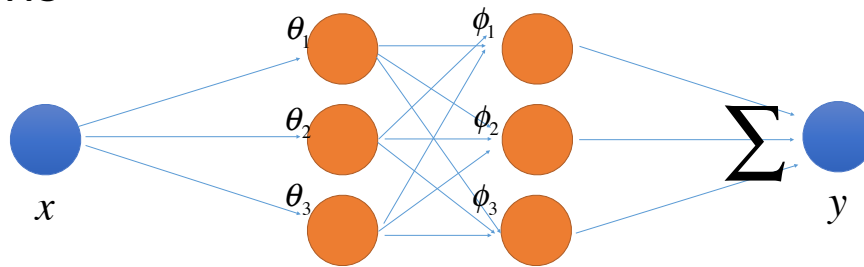
Assume $w_i = 1$ for all neurons, just learn the bias term b_i

Thin + Tall NN
With 6 neurons



Interpretation I: With the same number of neurons,
create combinatorial data flow

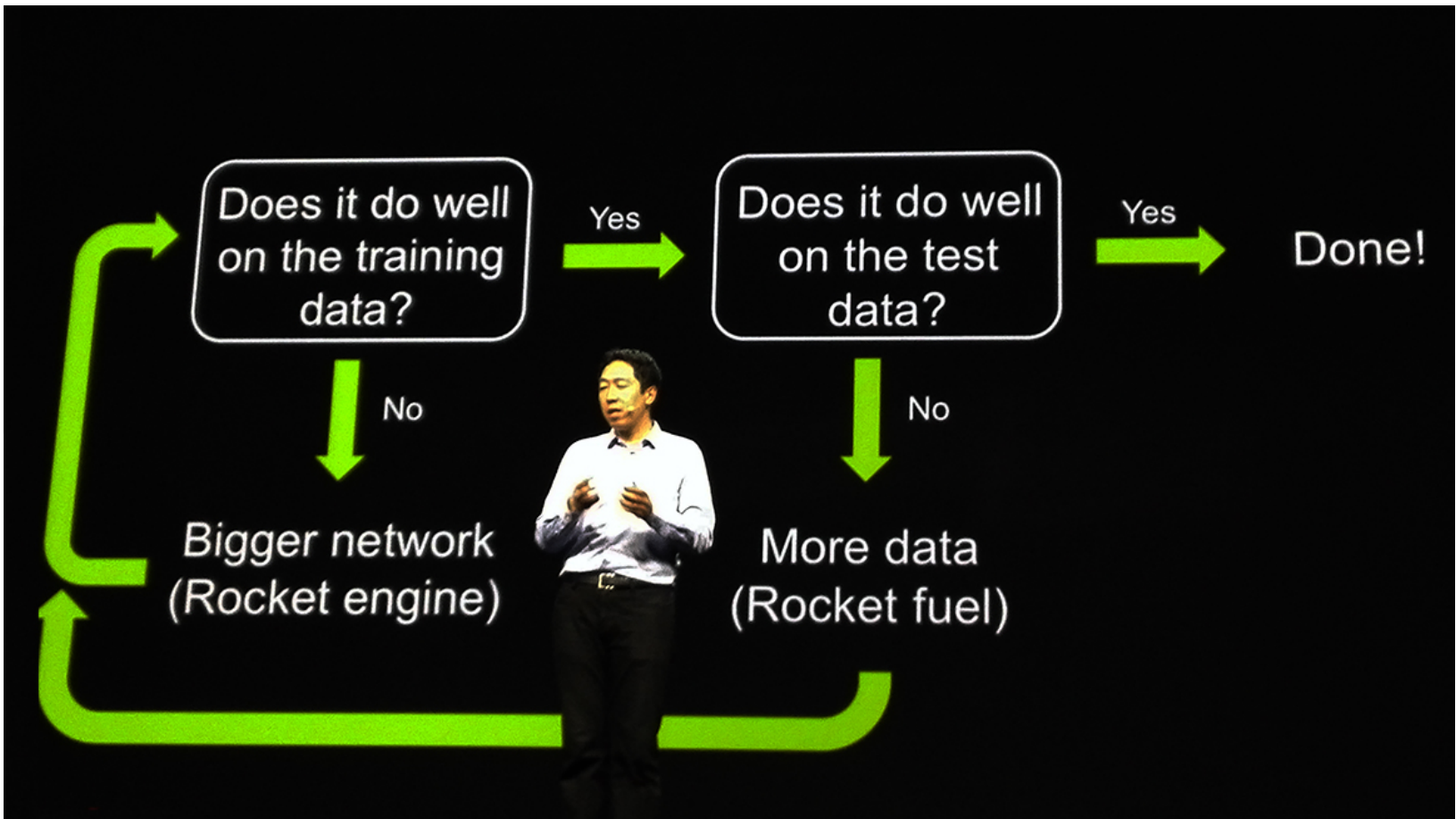
Thin + Tall NN
With 6 neurons



Interpretation I: With the same number of neurons,
create combinatorial data flow

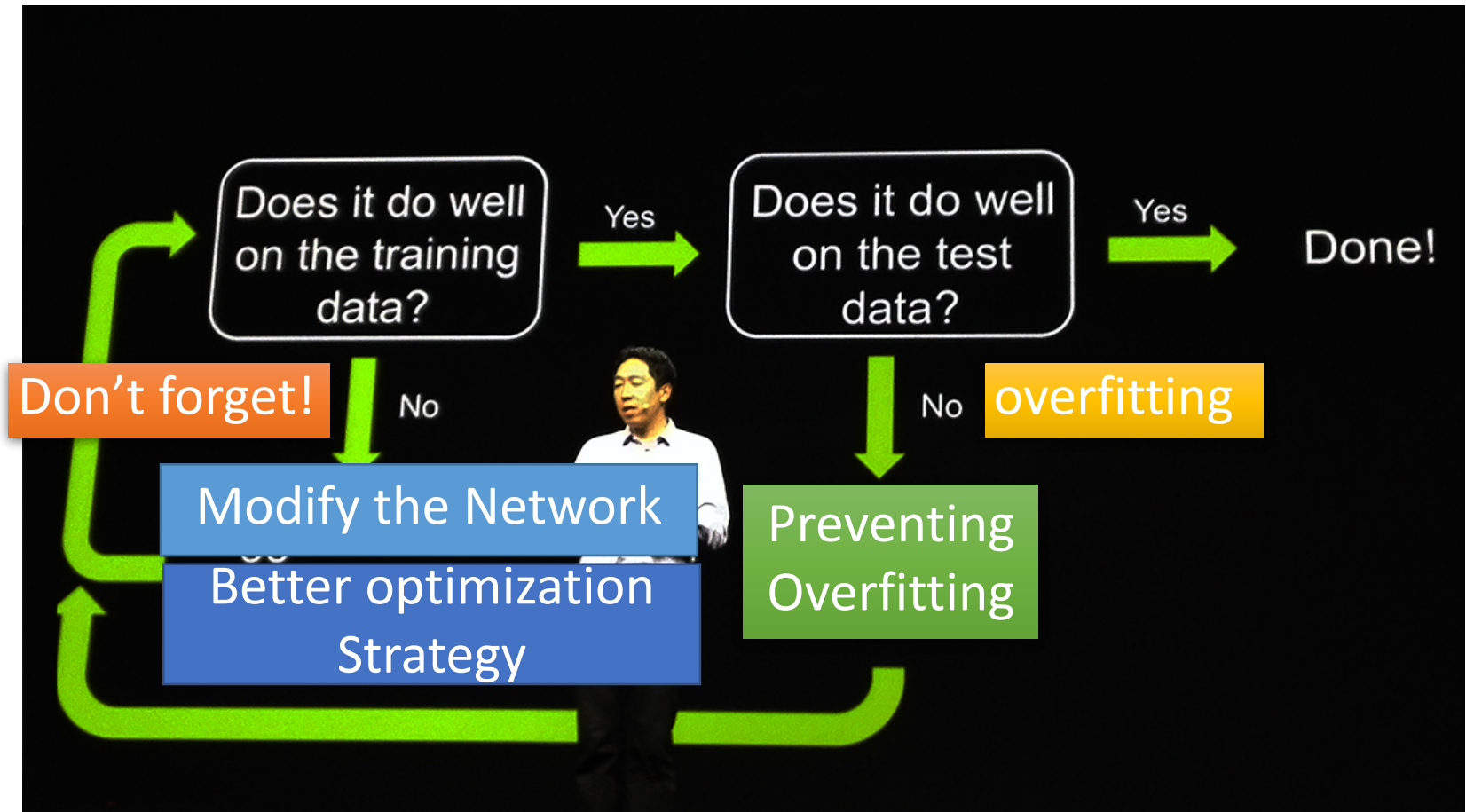
Interpretation II: Abstract data progressively
(edge-part-object)

Recipe for Learning



<http://www.gizmodo.com.au/2015/04/the-basic-recipe-for-machine-learning-explained-in-a-single-powerpoint-slide/>

Recipe for Learning



<http://www.gizmodo.com.au/2015/04/the-basic-recipe-for-machine-learning-explained-in-a-single-powerpoint-slide/>

Next lecture:

A big step-forward to reduce parameters of networks:

Convolutional Neural Network