# CS 466/566
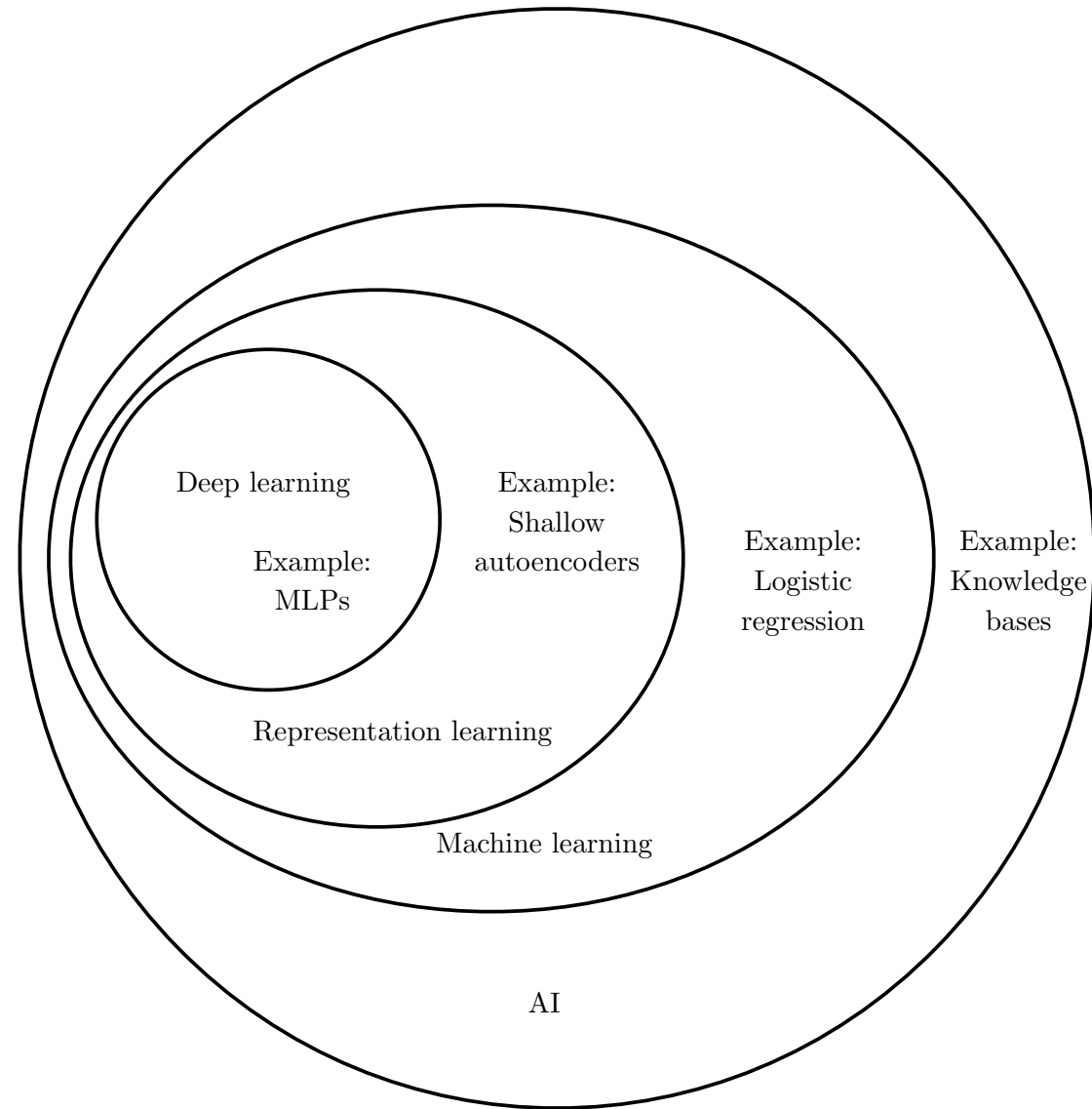# Introduction to Deep Learning

## Lecture 6 – Techniques for Training Better

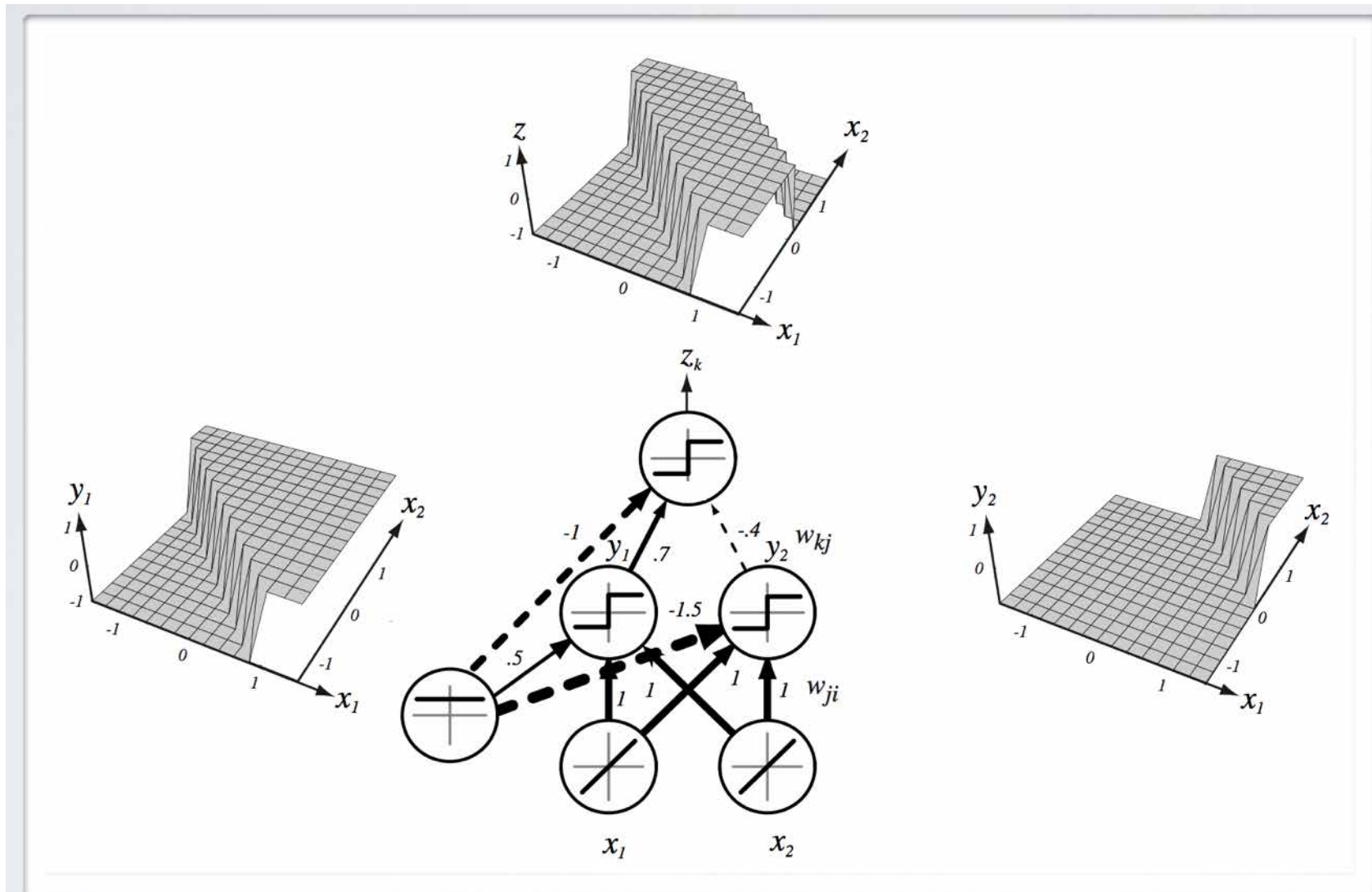Compiled from various resources

# Recall: ML and AI vs Deep learning



Deep learning

Example: MLPs

Example: Shallow autoencoders

Representation learning

Example: Logistic regression

Example: Knowledge bases

Machine learning

AI

# Recall: Capacity of a Neuron



range determined by $g(\cdot)$

bias $b$ only changes the position of the riff
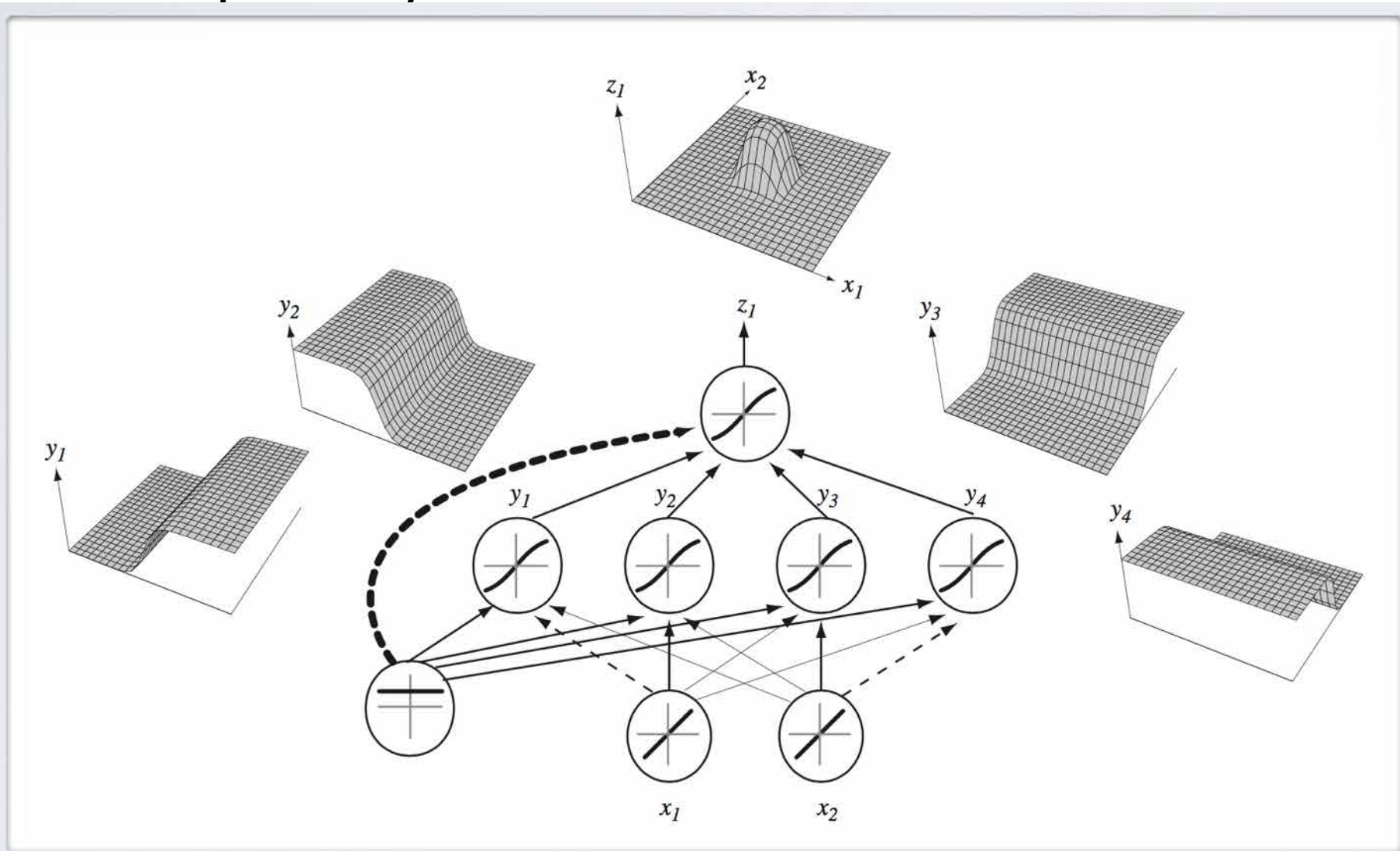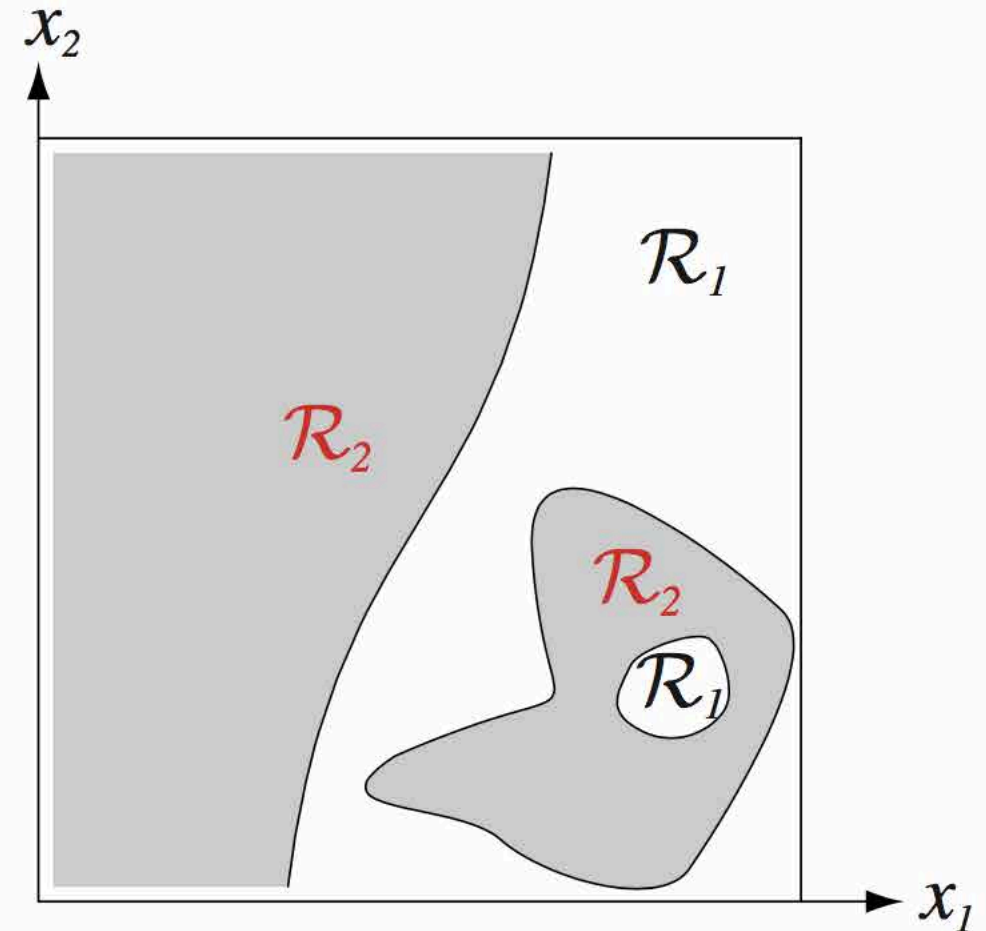
(from Pascal Vincent's slides)
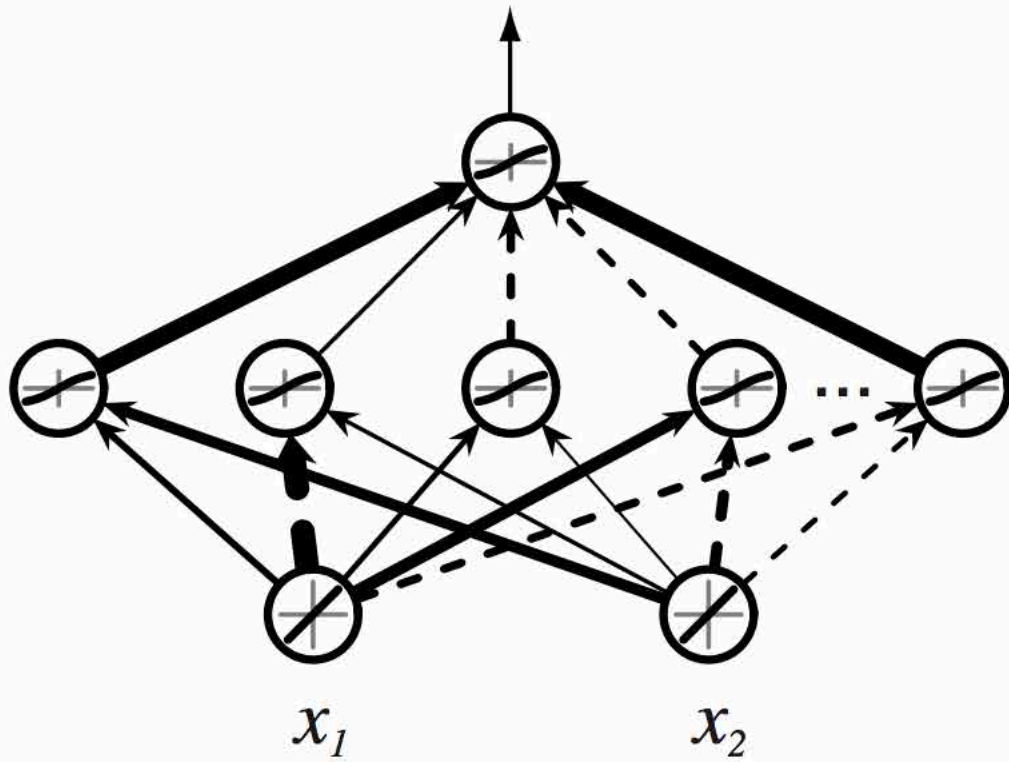
# Recall: Capacity of Neural Network



(from Pascal Vincent's slides)

# Recall: Capacity of Neural Network



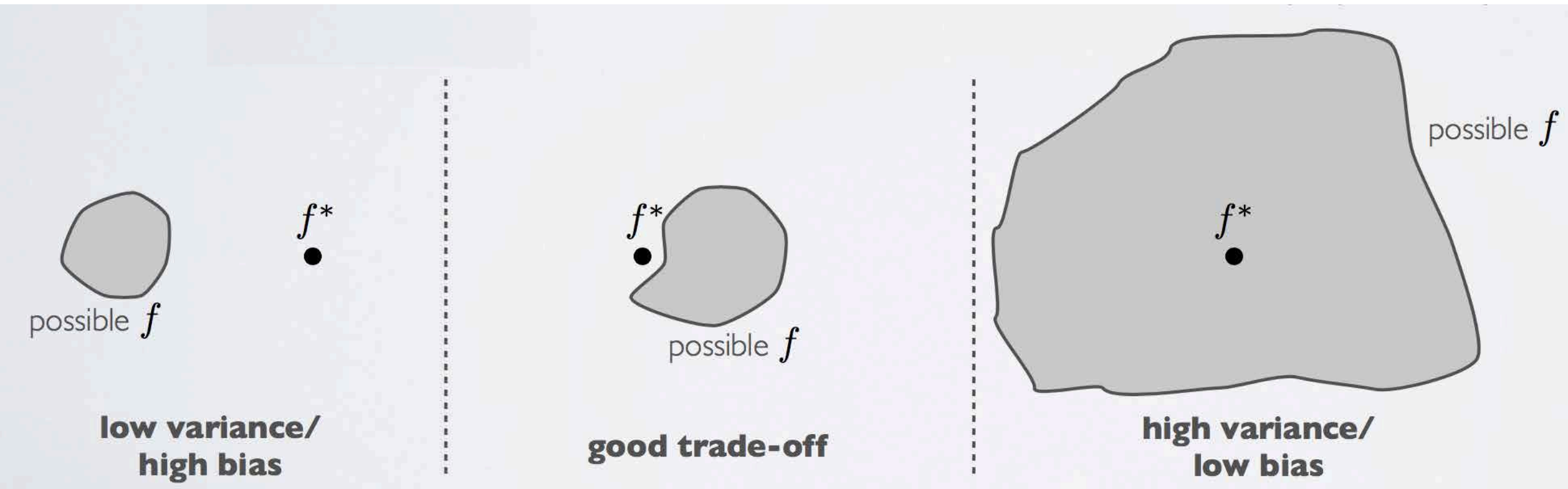(from Pascal Vincent's slides)

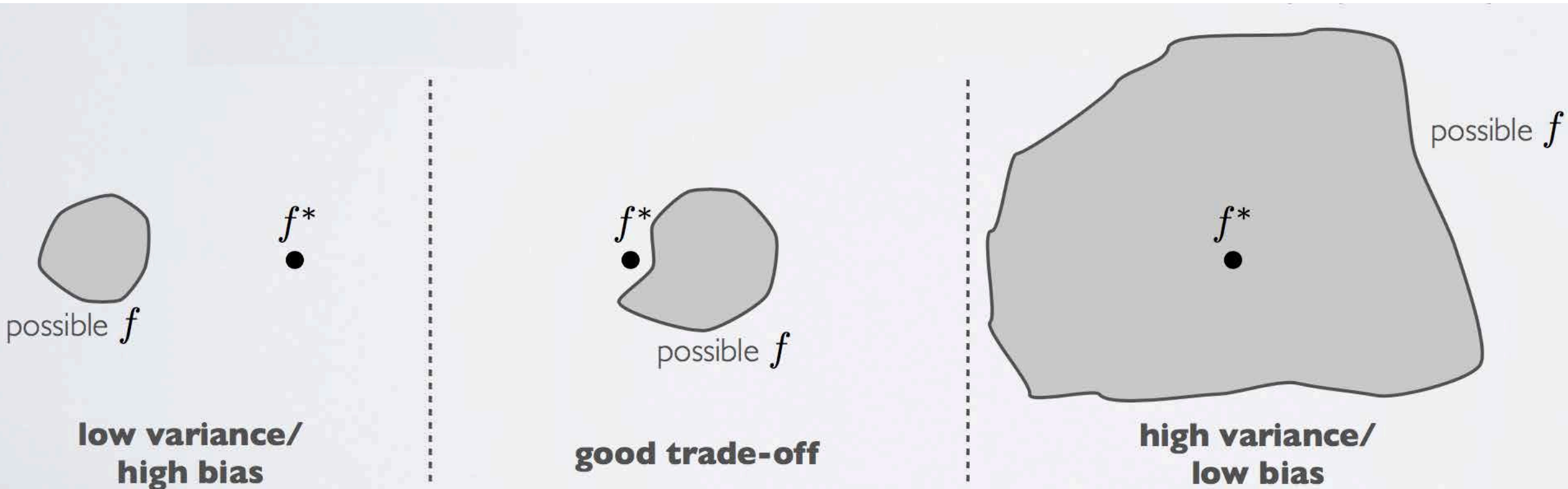# Recall: Capacity of Neural Network

# Bias – Variance Trade-off

- **Variance of trained model:** does it vary a lot if the training set changes?
- **Bias of trained model:** the average model close to the true solution
- **Generalization error:** (can be seen as) the sum of the (squared) bias and variance



low variance/
high bias

good trade-off

high variance/
low bias

# Training Problem: Overfitting

- **Variance of trained model:** does it vary a lot if the training set changes?

- **Bias of trained model:** the average model close to the true solution

- **Generalization error:** (can be seen as) the sum of the (squared) bias and variance

  - **Might be in high variance/low bias situation**



low variance/
high bias

good trade-off

high variance/
low bias

# Standard Neural Network Training Demo

- http://cs.stanford.edu/people/karpathy/svmjs/demo/demonn.html
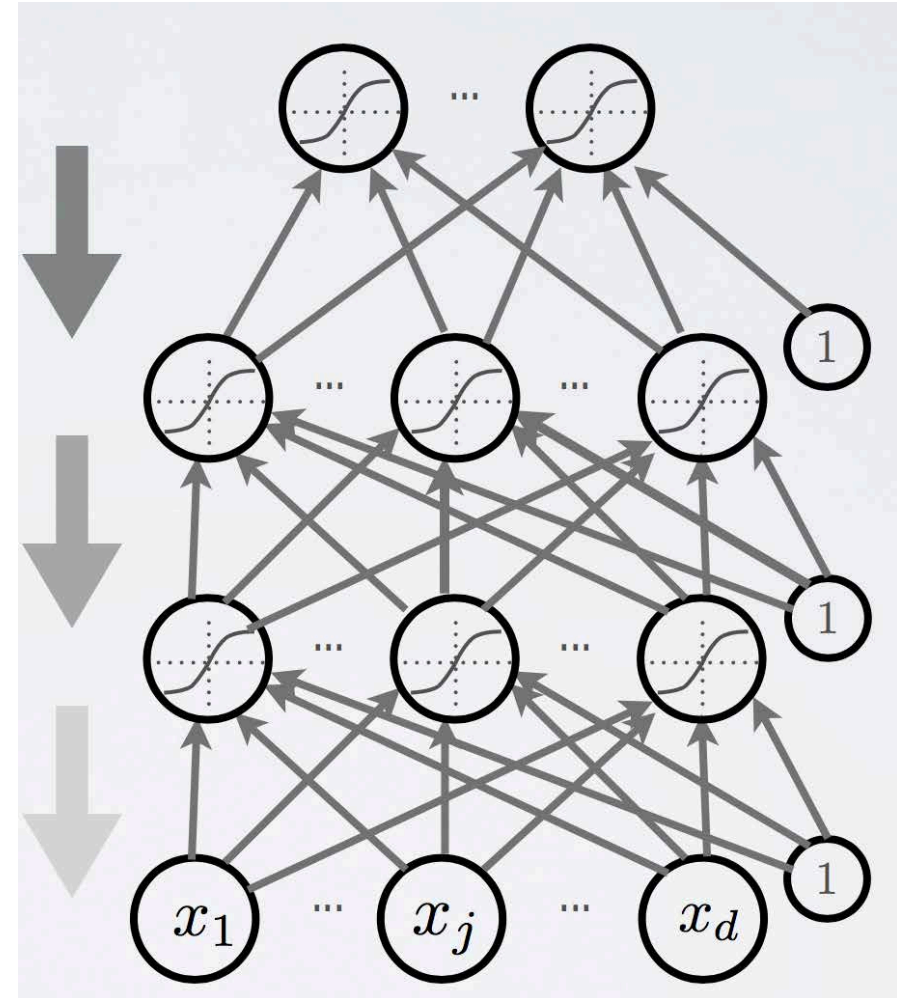
- We have lots of problems to tackle with during NN training

- Some of them are:
  - Vanishing gradient
  - Adversarial Samples
  - High variance / low bias networks (generally related to overfitting)

# Training Problem: Vanishing Gradient Problem

- Saturated units block gradient propagation.

- This is especially true for recurrent neural networks. Why?
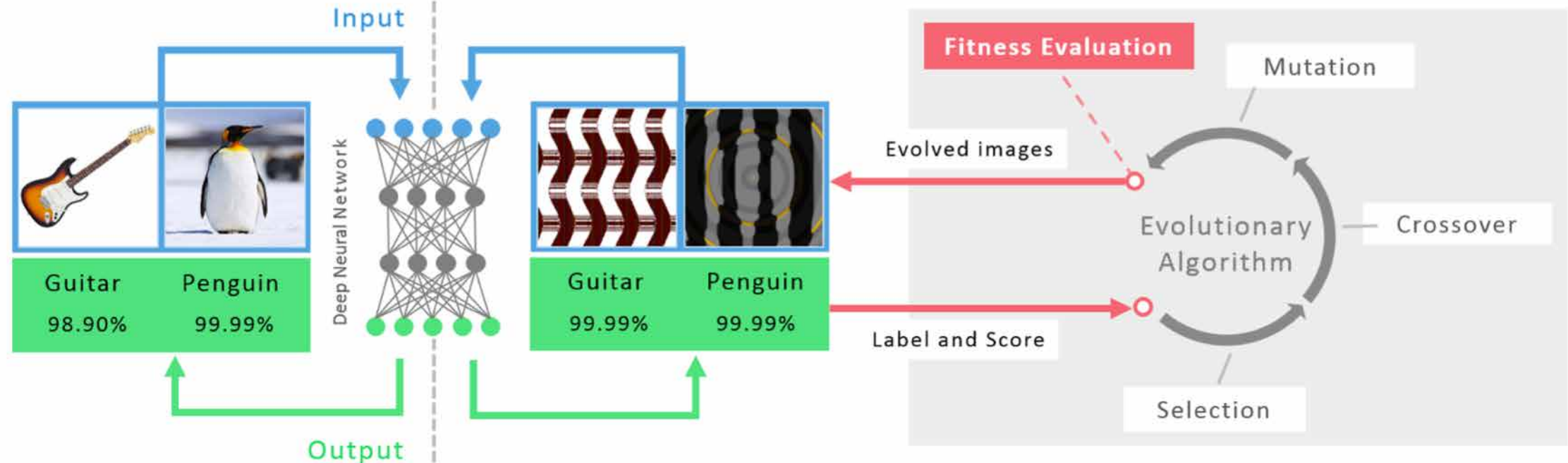
- Solution?

# Problem: Adversarial Examples

**1** State-of-the-art DNNs can recognize real images with high confidence
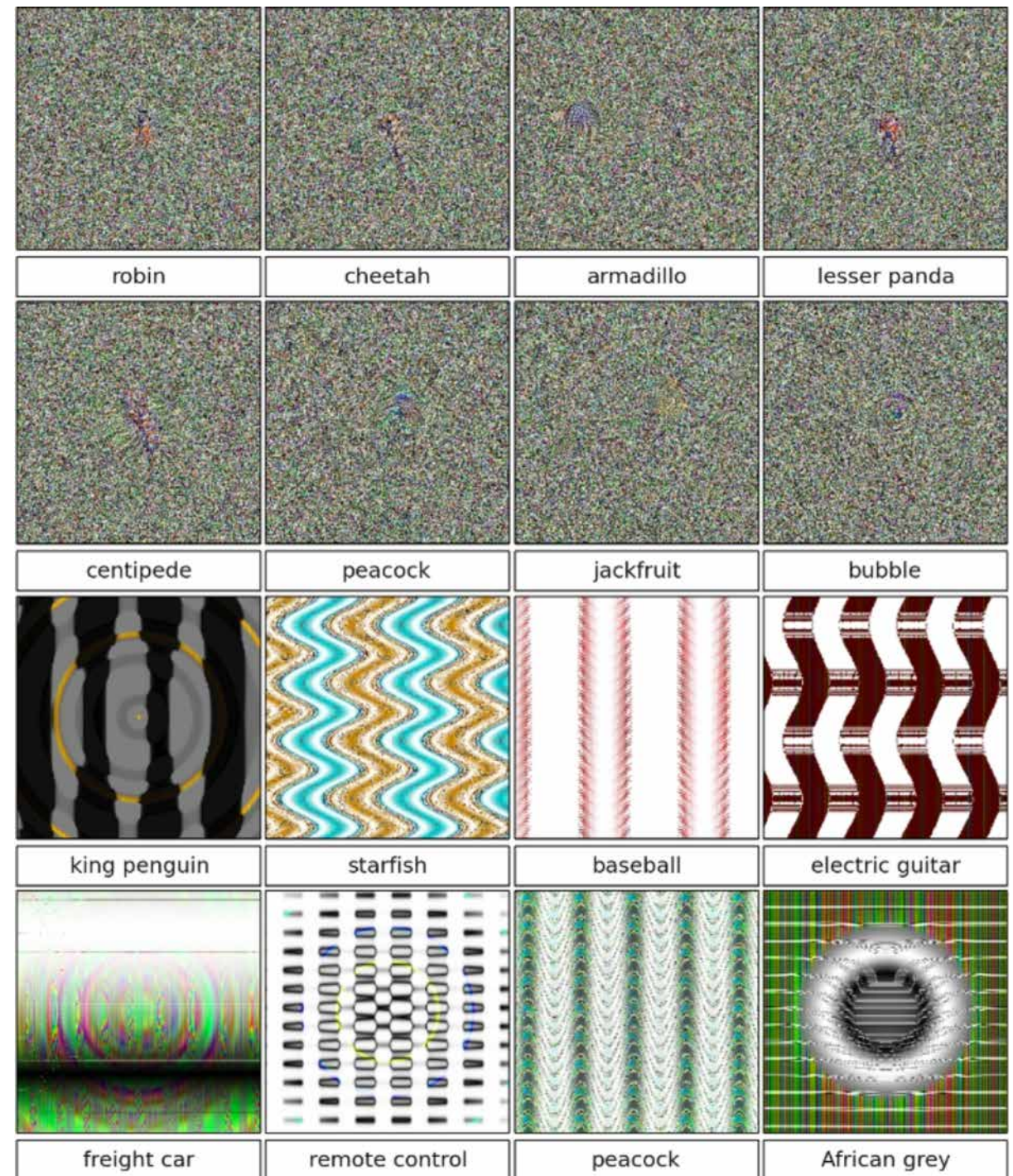
**2** But DNNs are also easily fooled: images can be produced that are unrecognizable to humans, but DNNs believe with 99.99% certainty are natural objects

Input

Deep Neural Network

Guitar 98.90%  Penguin 99.99%

Guitar 99.99%  Penguin 99.99%

Output

Fitness Evaluation

Mutation

Evolved images

Evolutionary Algorithm

Crossover

Label and Score

Selection

- Although state-of-the-art DNNs can increasingly recognize natural images, they also are easily fooled into declaring with near-certainty that unrecognizable images are familiar objects.

- Images that fool DNNs are produced by evolutionary algorithms that optimize images to generate high-confidence DNN predictions for each class in the dataset the DNN is trained on (here, ImageNet).

# Problem: Adversarial Examples

- Evolved images that are unrecognizable to humans but that state-of-the-art DNNs trained on ImageNet believe with ≥ 99.6% certainty to be a familiar object.

- This result highlights differences between how DNNs and humans recognize objects.

# Problem: Adversarial Examples



$+ .007 \times$

$=$

$\boldsymbol{x}$

$y =$"panda"
w/ 57.7%
confidence

$\mathrm{sign}(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$

"nematode"
w/ 8.2%
confidence

$\boldsymbol{x} +$
$\epsilon \, \mathrm{sign}(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$

"gibbon"
w/ 99.3 %
confidence

# Problem: Adversarial Examples in Real World



washer: 0.5398173

safe: 0.34602574
washer: 0.22088042

safe: 0.3719305
loudspeaker: 0.24184975

(a) Image from dataset  (b) Clean image  (c) Adv. image, $\epsilon = 4$  (d) Adv. image, $\epsilon = 8$

# Overfitting

- To overcome overfitting:
  - Pre-training (unsupervised learning): Auto-encoders, Stacked auto-encoders, Restricted Boltzmann Machines, etc...
  - We can use local receptive fields and shared weights, which makes the neural network a Convolutional Neural Network.
  - We can corrupt inputs (Denoising Auto-encoders)
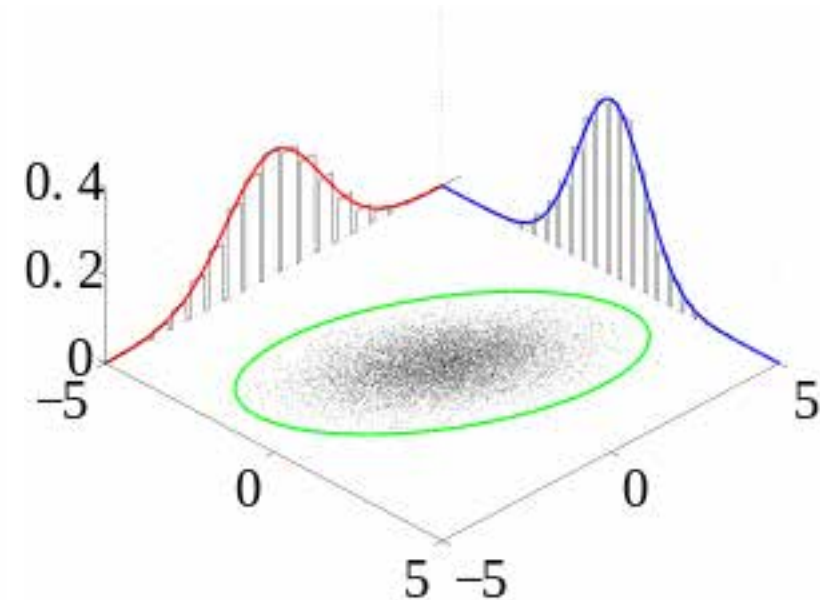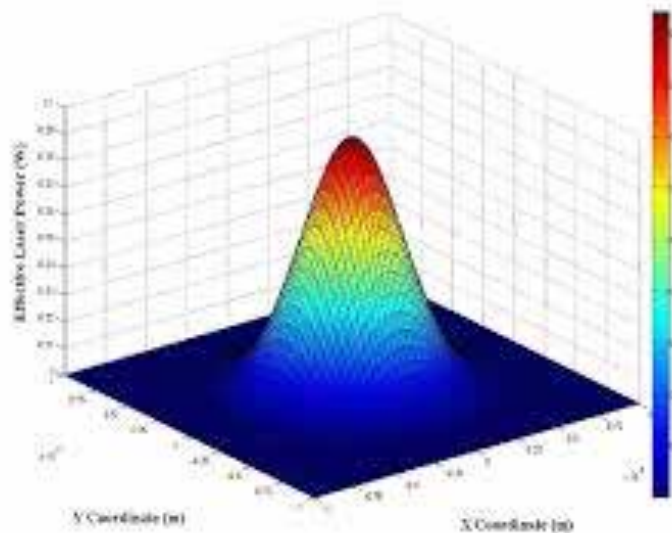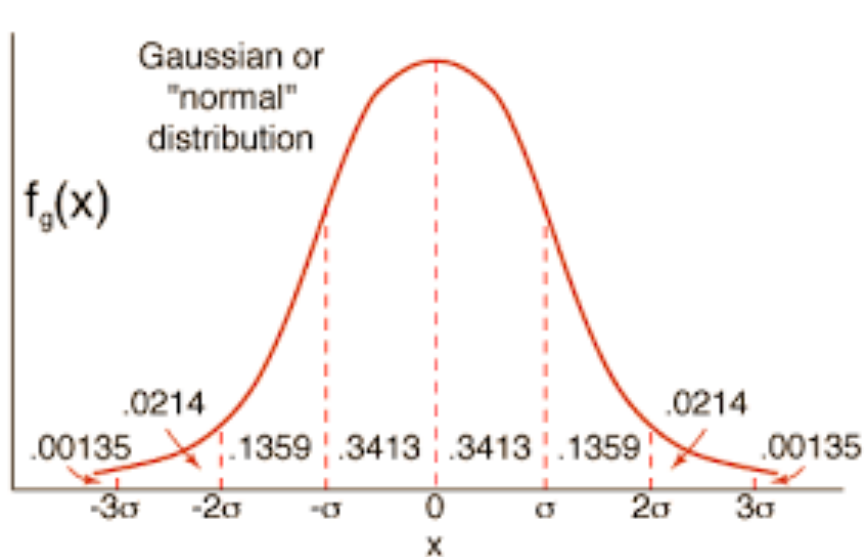  - We can use regularization (L1, L2, dropout, augmentation)

# Regularization

- Regularization is any modification we make to a learning algorithm that is intended to reduce its generalization error but not its training error.


- L2 regularization

- L1 regularization

- Dropout

- Data augmentation

# Regularization: L2 Regularization

- Penalizes the gradient term in back-prop algorithm:

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

- Gradient: $\nabla_{\mathbf{W}^{(k)}} \Omega(\boldsymbol{\theta}) = 2\mathbf{W}^{(k)}$

  - Only applied on weights, not biases
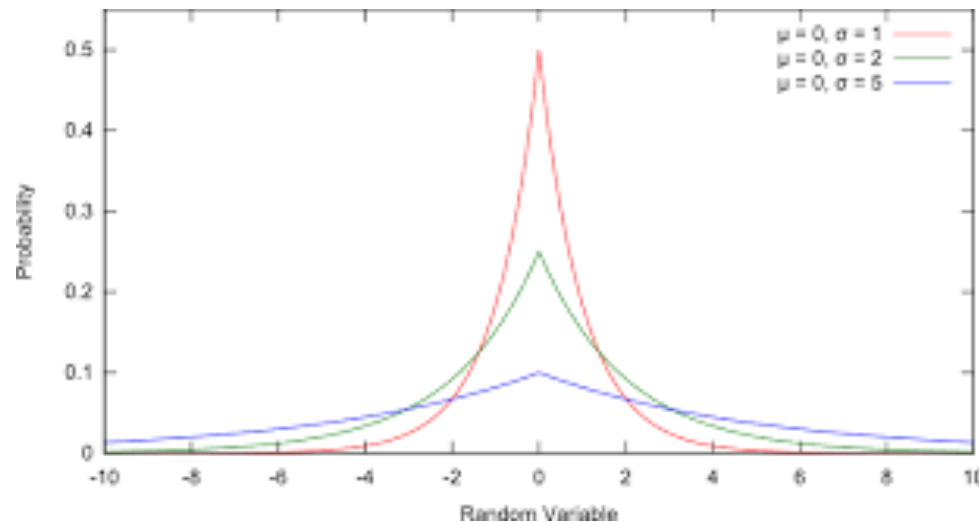  - Can be interpreted as having a Gaussian Prior over weight values.

# Regularization: L1 Regularization

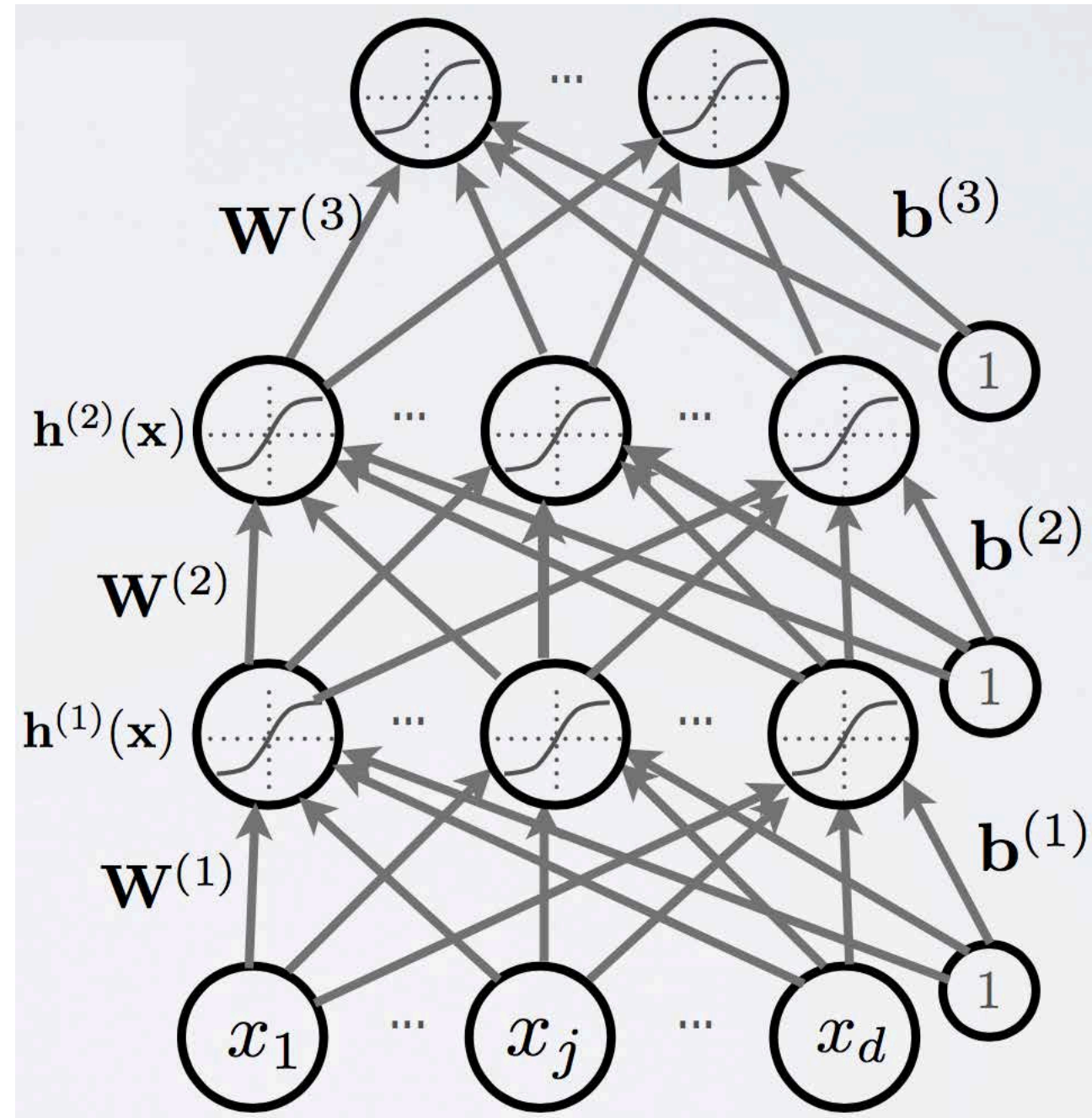- Penalizes the gradient term in back-prop algorithm:

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

- Gradient: $\nabla_{\mathbf{W}^{(k)}} \Omega(\boldsymbol{\theta}) = \text{sign}(\mathbf{W}^{(k)})$
  - Only applied on weights, not biases (again)
  - Unlike L2, L1 regularization will make some of the weights exactly 0.
  - Can be interpreted as having a Laplacian Prior over weight values.

# Regularization: Dropout

- Idea: Cripple neural network by removing hidden units randomly.
  - Each hidden unit's activation is set to 0 with a probability, independently.
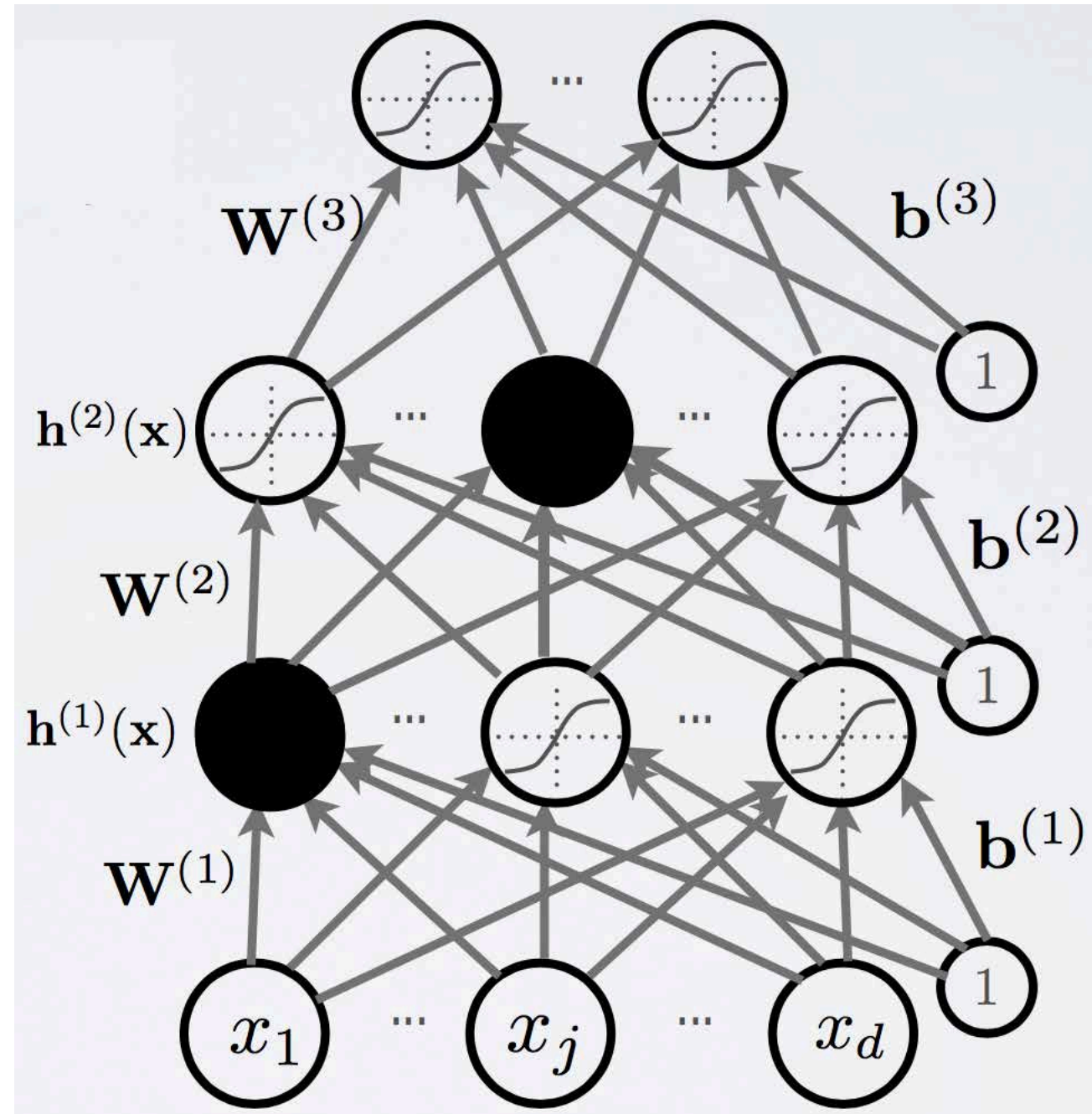  - Usually 0.5 works well.

# Regularization: Dropout

- Idea: Cripple neural network by removing hidden units randomly.
  - Each hidden unit's activation is set to 0 with a probability, independently.
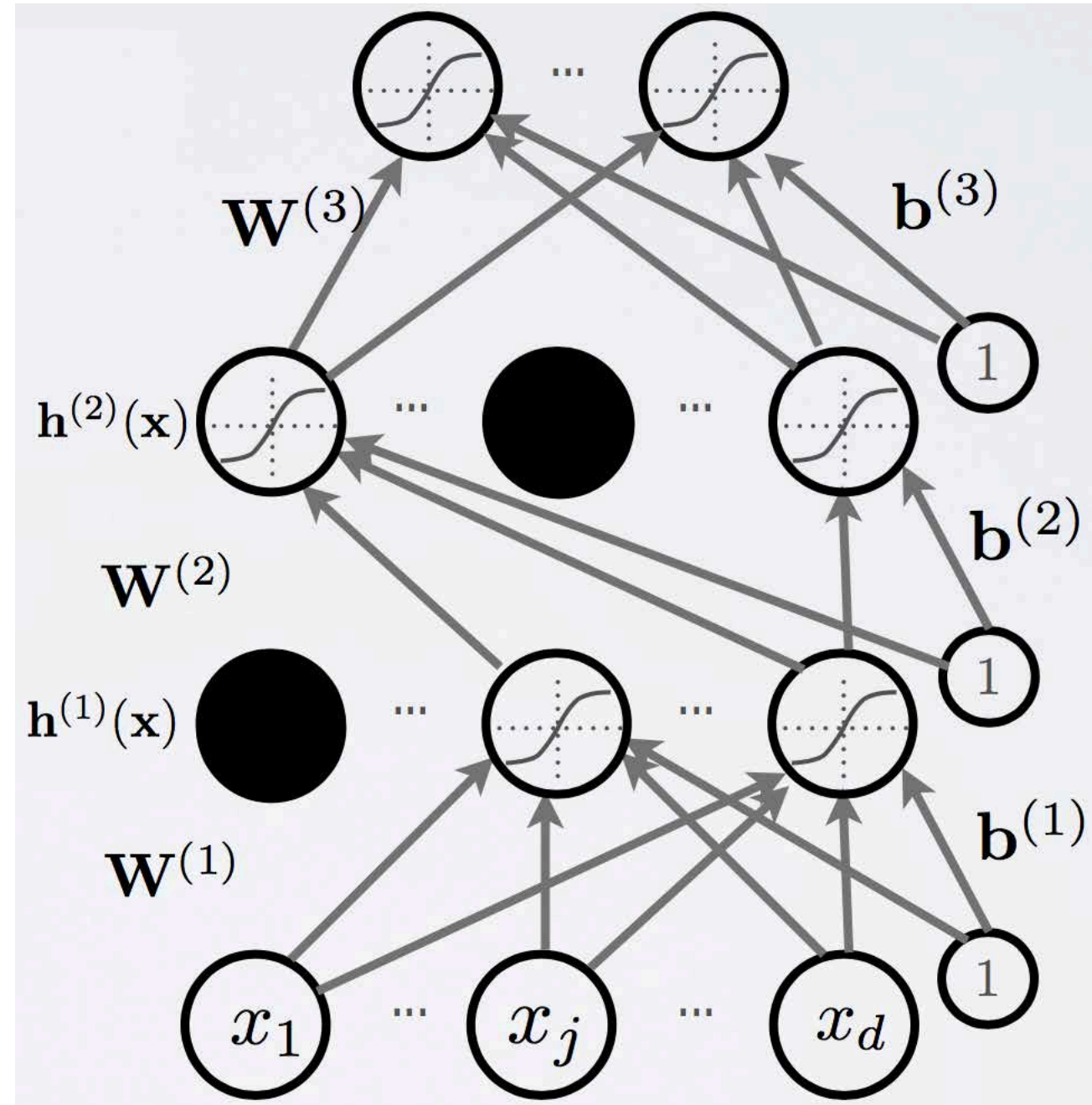  - Usually 0.5 works well.

# Regularization: Dropout

- Idea: Cripple neural network by removing hidden units randomly.
  - Hidden units cannot cooperate with each other in a specific layer.
  - Therefore, hidden units must be more generally useful.

# Regularization: Dropout

- At test time, we replace the masks by their expectation
  - this is simply the constant vector 0.5 if dropout probability is 0.5
  - for single hidden layer, can show this is 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

# Regularization: Dropout



Base network

Ensemble of subnetworks

# Regularization: Dataset Augmentation

# Possible Variations on MNIST

# Impact of Pre-training
## (intentional underfitting scenario)

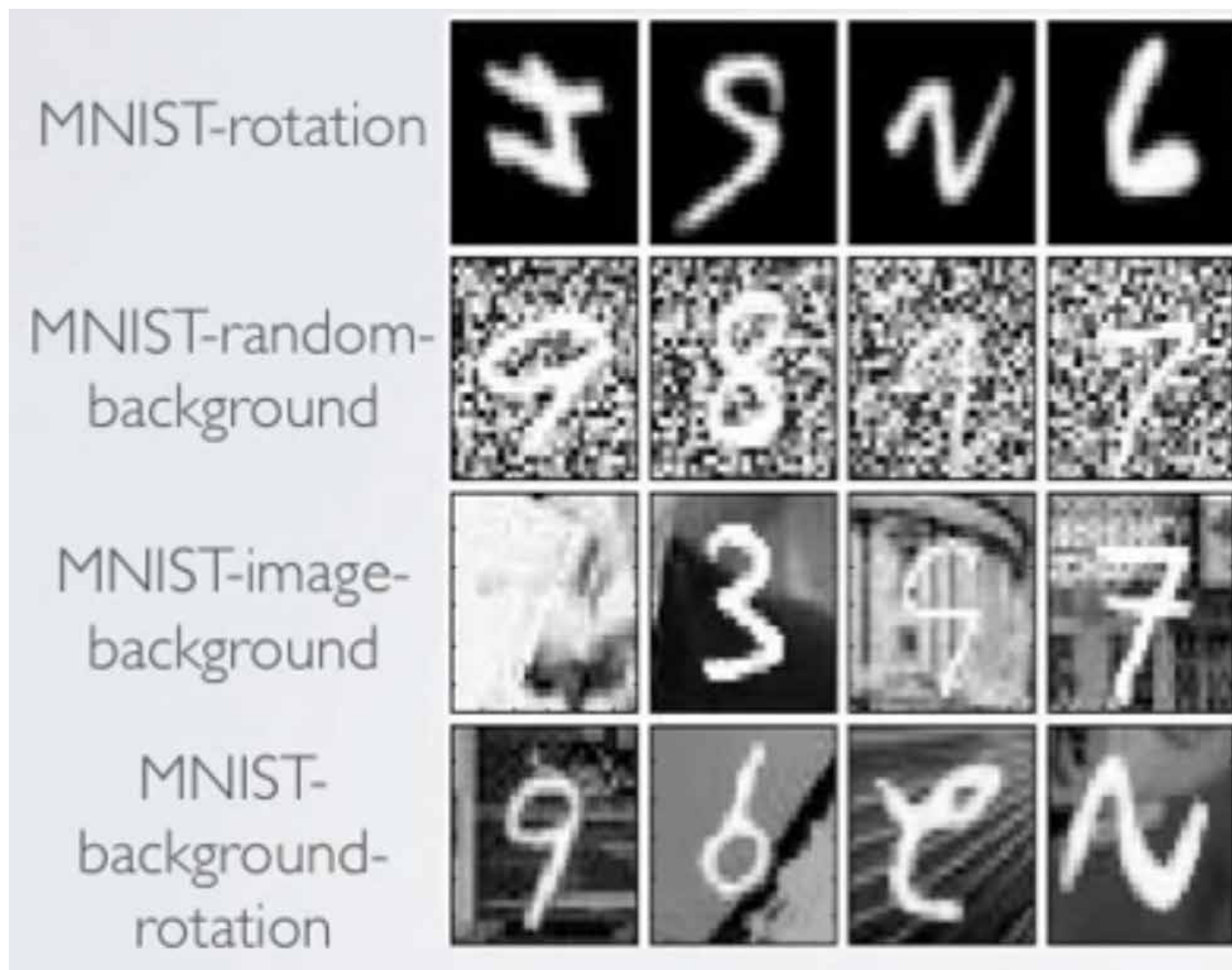| Network | | MNIST-small | MNIST-rotation |
|---|---|---|---|
| Type | Depth | classif. test error | classif. test error |
| Deep net | 1 | **4.14** % $\pm 0.17$ | 15.22 % $\pm 0.31$ |
| | 2 | **4.03** % $\pm 0.17$ | **10.63** % $\pm 0.27$ |
| | 3 | **4.24** % $\pm 0.18$ | 11.98 % $\pm 0.28$ |
| | 4 | 4.47 % $\pm 0.18$ | 11.73 % $\pm 0.29$ |
| Deep net + autoencoder | 1 | 3.87 % $\pm 0.17$ | 11.43% $\pm 0.28$ |
| | 2 | **3.38** % $\pm 0.16$ | 9.88 % $\pm 0.26$ |
| | 3 | **3.37** % $\pm 0.16$ | **9.22** % $\pm 0.25$ |
| | 4 | **3.39** % $\pm 0.16$ | **9.20** % $\pm 0.25$ |
| Deep net + RBM | 1 | 3.17 % $\pm 0.15$ | 10.47 % $\pm 0.27$ |
| | 2 | **2.74** % $\pm 0.14$ | 9.54 % $\pm 0.26$ |
| | 3 | **2.71** % $\pm 0.14$ | **8.80** % $\pm 0.25$ |
| | 4 | **2.72** % $\pm 0.14$ | **8.83** % $\pm 0.24$ |

# Performance on Different Datasets

| Dataset | $\mathbf{SVM}_{rbf}$ | Stacked Autoencoders $\mathbf{SAA\text{-}3}$ | Stacked RBMS $\mathbf{DBN\text{-}3}$ | Stacked Denoising Autoencoders $\mathbf{SdA\text{-}3}\ (\nu)$ |
|---|---|---|---|---|
| *basic* | $\mathbf{3.03{\pm}0.15}$ | $3.46{\pm}0.16$ | $3.11{\pm}0.15$ | $\mathbf{2.80{\pm}0.14}$ (10%) |
| *rot* | $11.11{\pm}0.28$ | $\mathbf{10.30{\pm}0.27}$ | $\mathbf{10.30{\pm}0.27}$ | $\mathbf{10.29{\pm}0.27}$ (10%) |
| *bg-rand* | $14.58{\pm}0.31$ | $11.28{\pm}0.28$ | $\mathbf{6.73{\pm}0.22}$ | $10.38{\pm}0.27$ (40%) |
| *bg-img* | $22.61{\pm}0.37$ | $23.00{\pm}0.37$ | $\mathbf{16.31{\pm}0.32}$ | $\mathbf{16.68{\pm}0.33}$ (25%) |
| *rot-bg-img* | $55.18{\pm}0.44$ | $51.93{\pm}0.44$ | $47.39{\pm}0.44$ | $\mathbf{44.49{\pm}0.44}$ (25%) |
| *rect* | $\mathbf{2.15{\pm}0.13}$ | $2.41{\pm}0.13$ | $2.60{\pm}0.14$ | $\mathbf{1.99{\pm}0.12}$ (10%) |
| *rect-img* | $24.04{\pm}0.37$ | $24.05{\pm}0.37$ | $22.50{\pm}0.37$ | $\mathbf{21.59{\pm}0.36}$ (25%) |
| *convex* | $19.13{\pm}0.34$ | $\mathbf{18.41{\pm}0.34}$ | $\mathbf{18.63{\pm}0.34}$ | $\mathbf{19.06{\pm}0.34}$ (10%) |

Extracting and Composing Robust Features with Denoising Autoencoders, Vincent, Larochelle, Bengio and Manzagol, 2008.