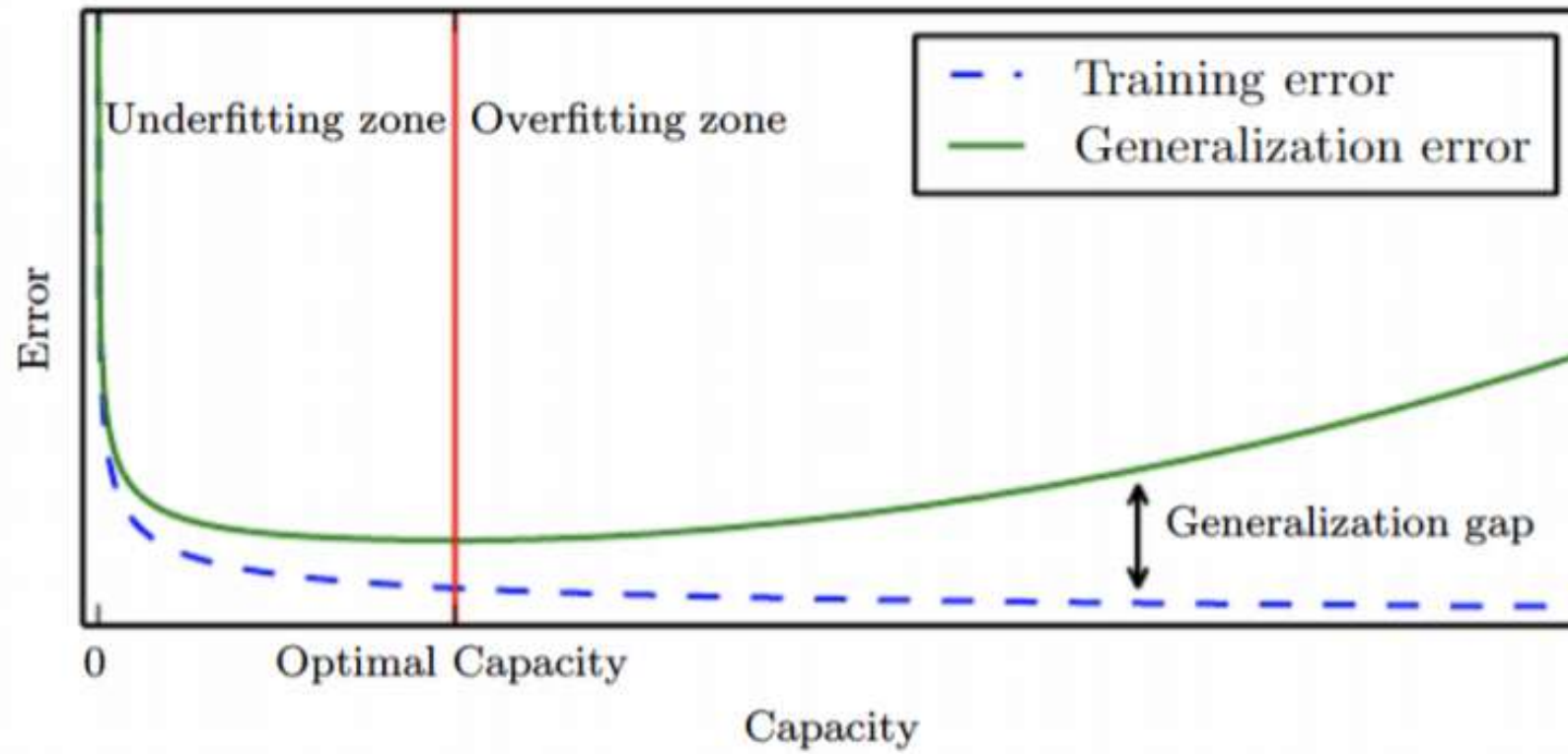


# Regularization

Deep Learning

# Why Regularization?



# Definition: Regularization

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

# Model Types and Regularization

- Three types of model families
  1. Excludes the true data generating process
    - Implies underfitting and inducing high bias
  2. Matches the true data generating process
  3. Overfits
    - Includes true data generating process but also many other processes
- Goal of regularization is to take model from third regime to second

# Norm Penalty

- When our training algorithm minimizes the regularized objective function

$$\tilde{J}(\theta; X, y) = J(\theta; X, y) + \alpha\Omega(\theta)$$

- it will decrease both the original objective  $J$  on the training data and some measure of the size of the parameters  $\theta$
- Different choices of the parameter norm  $\Omega$  can result in different solutions preferred
  - We discuss effects of various norms

## $L^2$ parameter Regularization

- Simplest and most common kind
- Called *Weight decay*
- Drives weights closer to the origin
  - by adding a regularization term to the objective function
- In other communities also known as *ridge regression* or *Tikhonov regularization*

$$\Omega(\theta) = \frac{1}{2} ||w||_2^2$$



# Gradient of Regularized Objective

- Objective function (with no bias parameter)

$$\tilde{J}(w; X, y) = \frac{\alpha}{2} w^T w + J(w; X, y)$$

- Corresponding parameter gradient

$$\nabla_w \tilde{J}(w; X, y) = \alpha w + \nabla_w J(w; X, y)$$

- To perform single gradient step, perform update:

$$w \leftarrow w - \varepsilon \left( \alpha w + \nabla_w J(w; X, y) \right)$$

- Written another way, the update is

$$w \leftarrow (1 - \varepsilon \alpha) w - \varepsilon \nabla_w J(w; X, y)$$

- We have modified learning rule to shrink  $w$  by constant factor  $1 - \varepsilon \alpha$  at each step

# Effect of $L^2$ regularization on optimal $w$

Objective function:

$$\tilde{J}(w; X, y) = \frac{\alpha}{2} w^T w + J(w; X, y)$$

Solid ellipses:

contours of equal value of  
unregularized objective

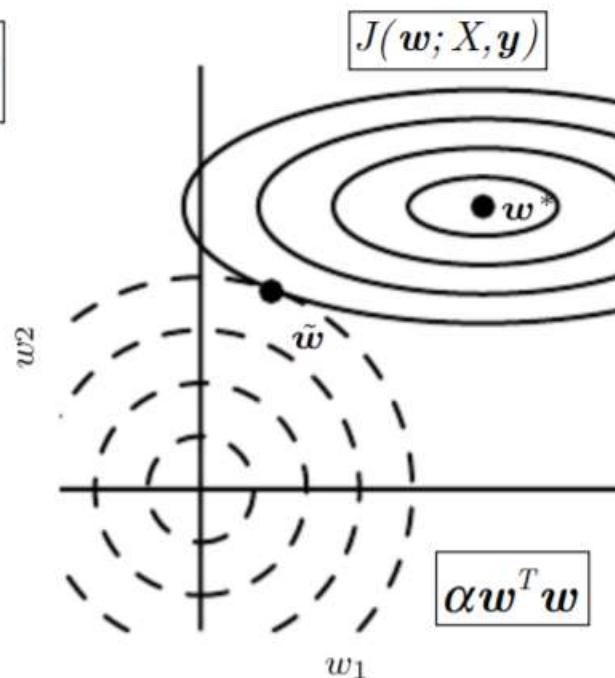
$$J(w; X, y)$$

Dotted circles:

contours of equal value of  $L^2$   
regularizer

$$\alpha w^T w$$

At point  $\tilde{w}$  competing objectives  
reach equilibrium



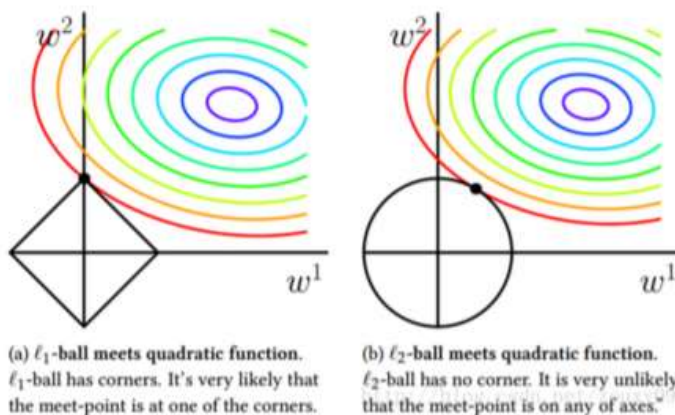


# $L^1$ Regularization

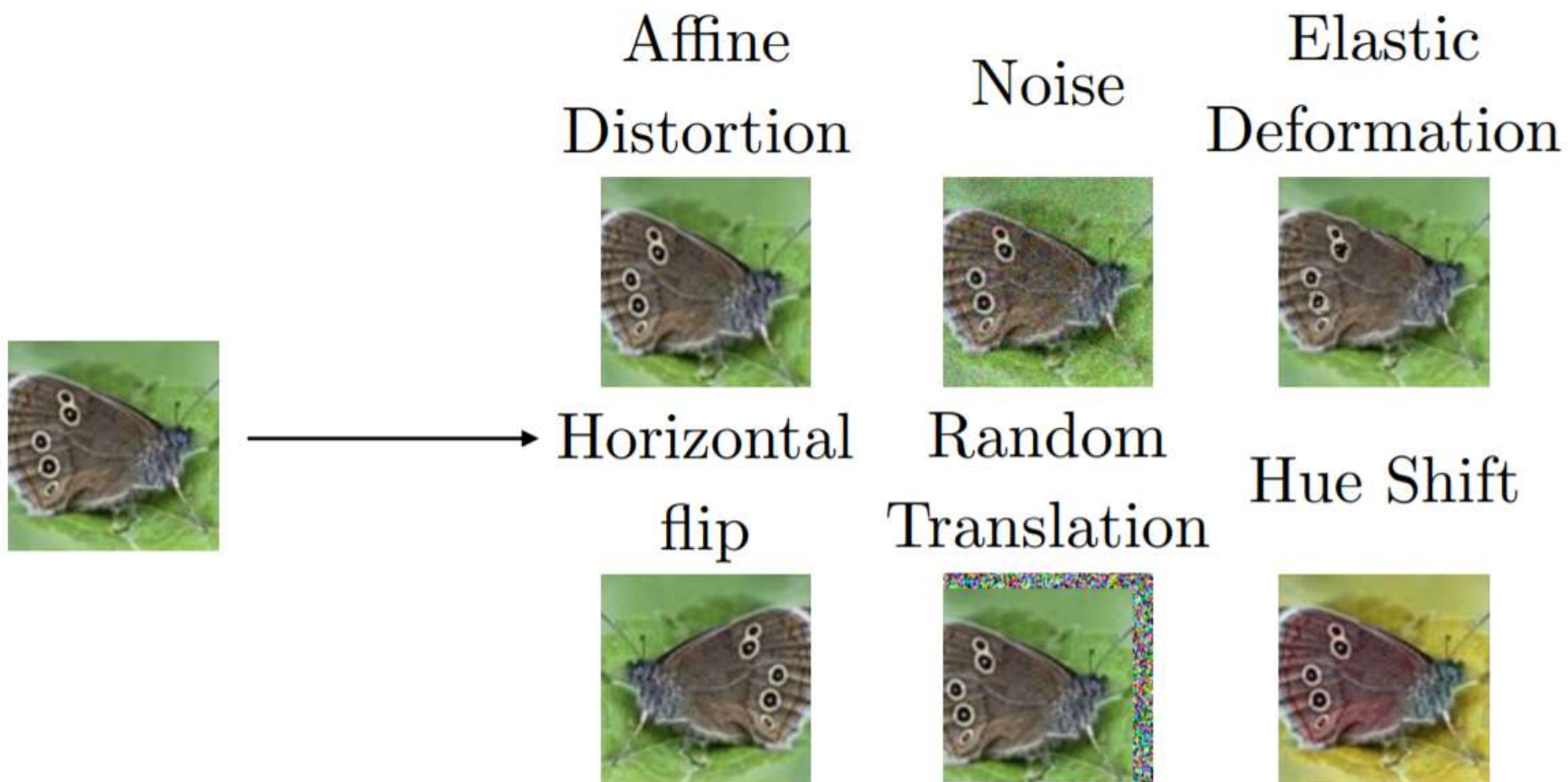
- $L^2$  weight decay is common weight decay
- Other ways to penalize model parameter size
- $L^1$  regularization is defined as

$$\Omega(\theta) = \|w\|_1 = \sum_i |w_i|$$

– which sums the absolute values of parameters



# Data Augmentation

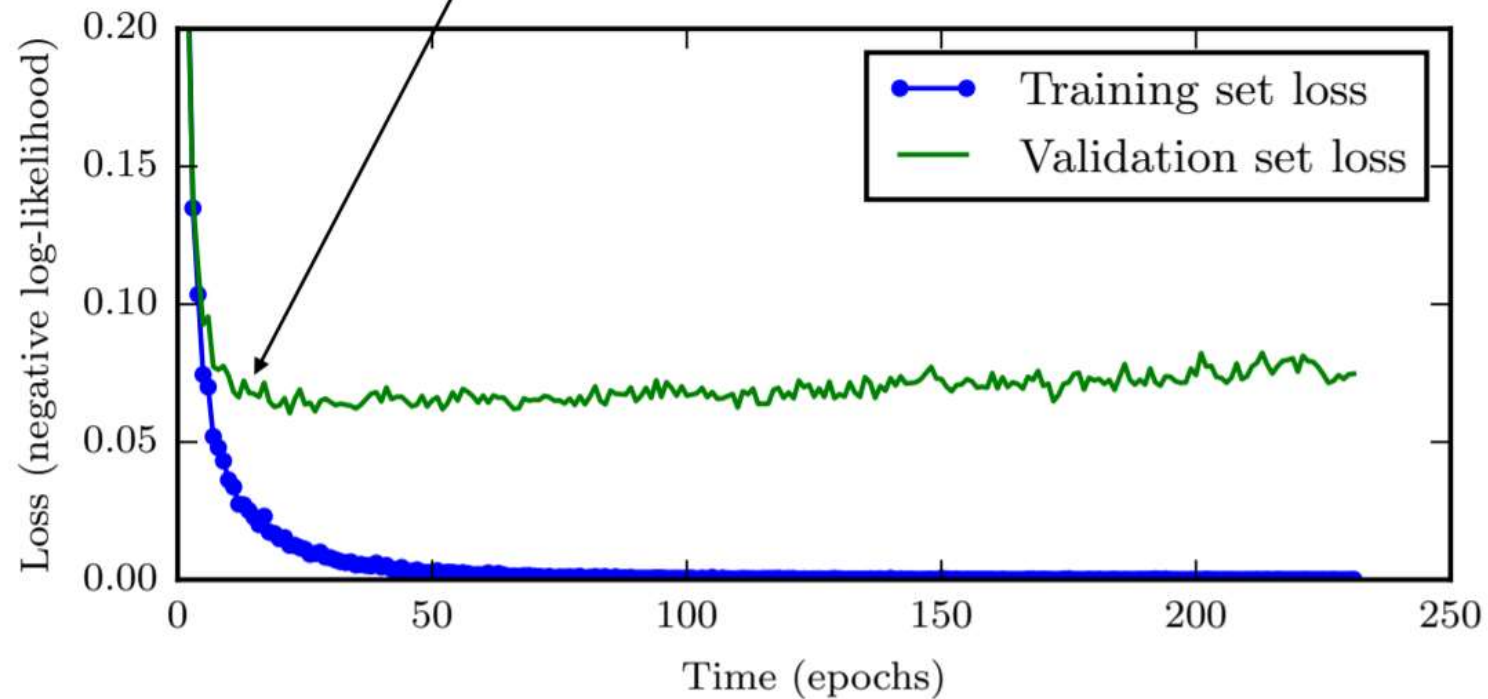


# Caution in Data Augmentation

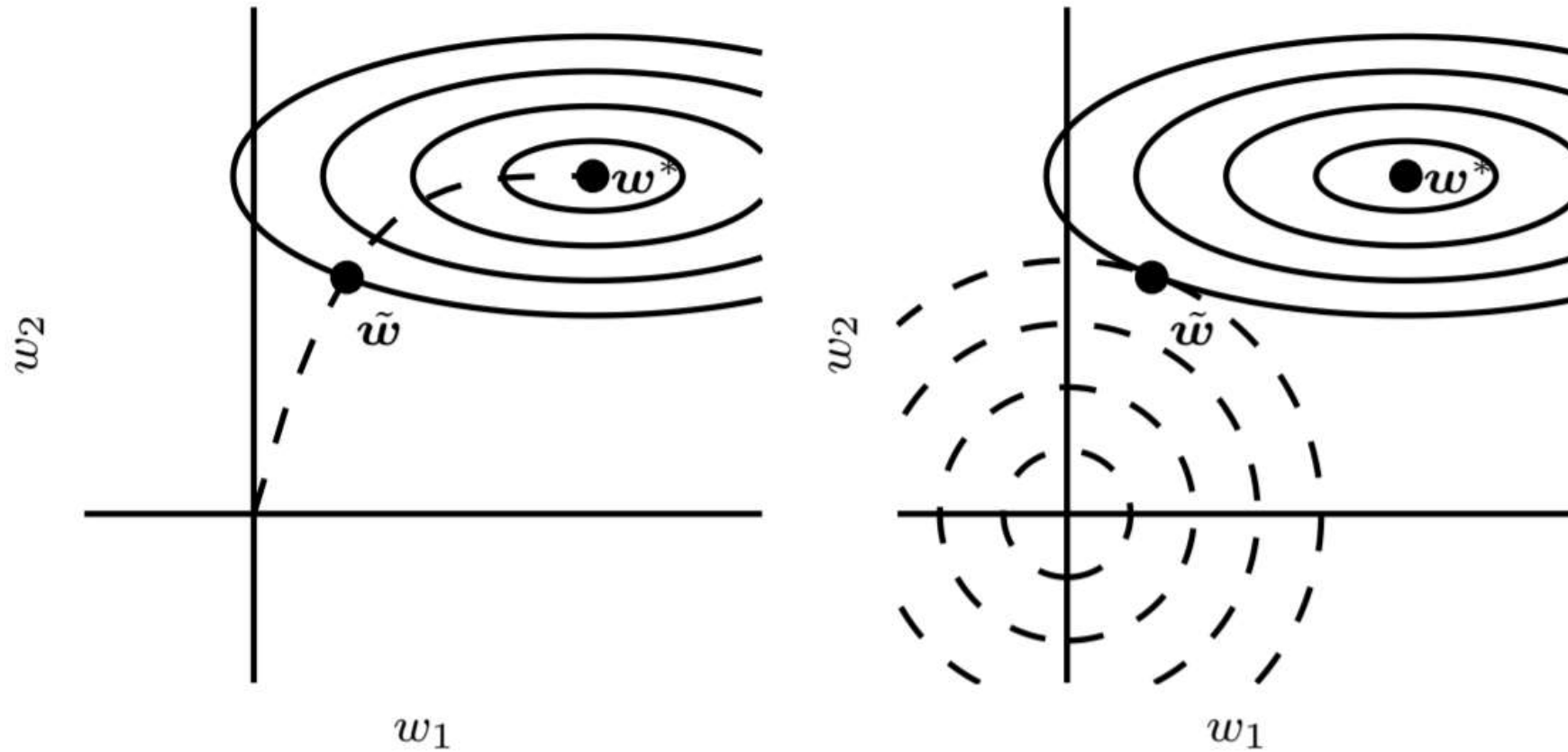
- Not apply transformation that would change the class
- OCR example: 'b' vs 'd' and '6' vs '9'
  - Horizontal flips and 180 degree rotations are not appropriate ways
- Some transformations are not easy to perform
  - Out of plane rotation cannot be implemented as a simple geometric operation on pixels

# Early Stopping

Early stopping: terminate while validation set performance is better



# Early Stopping and Weight Decay



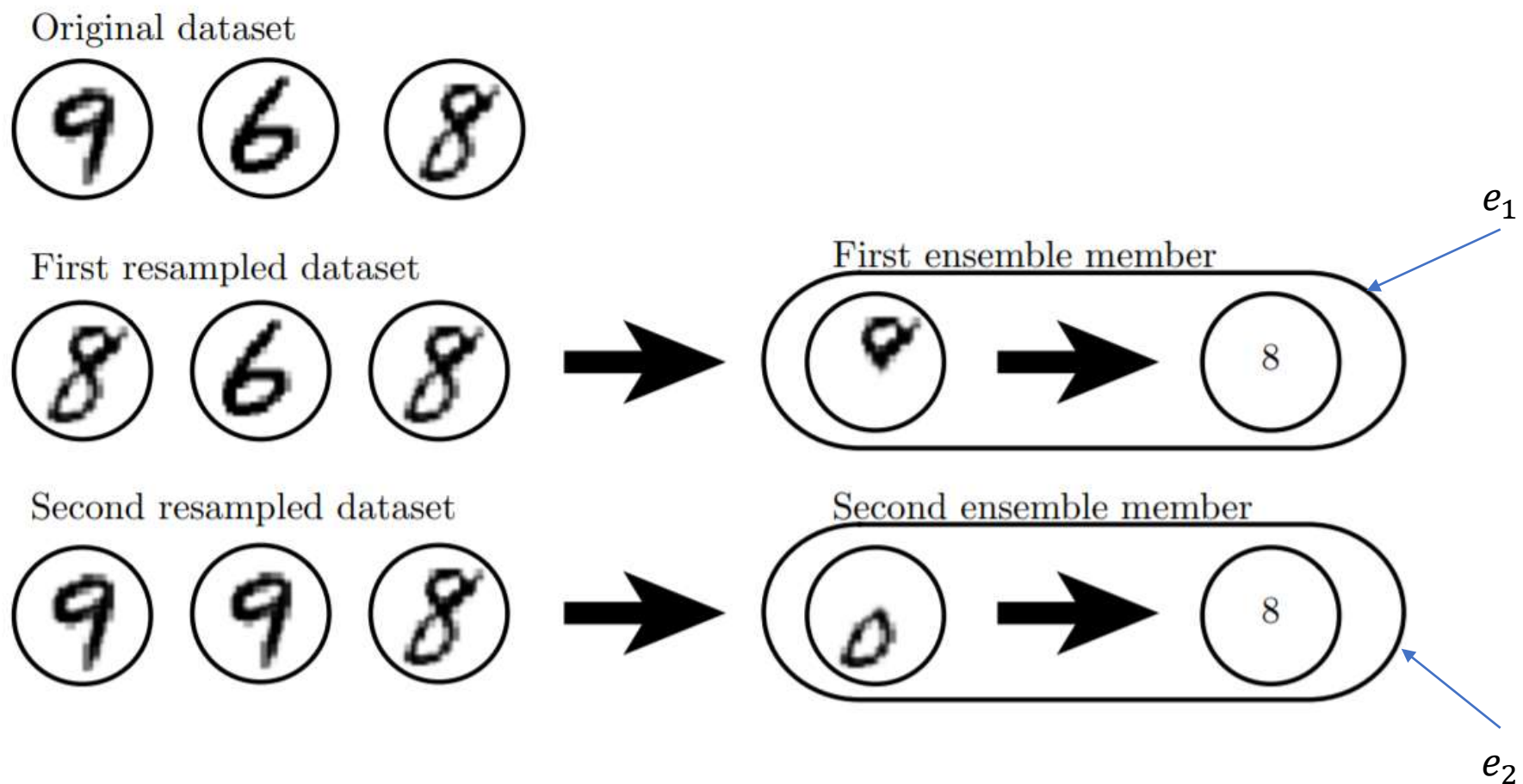


# Bagging/Model Ensemble

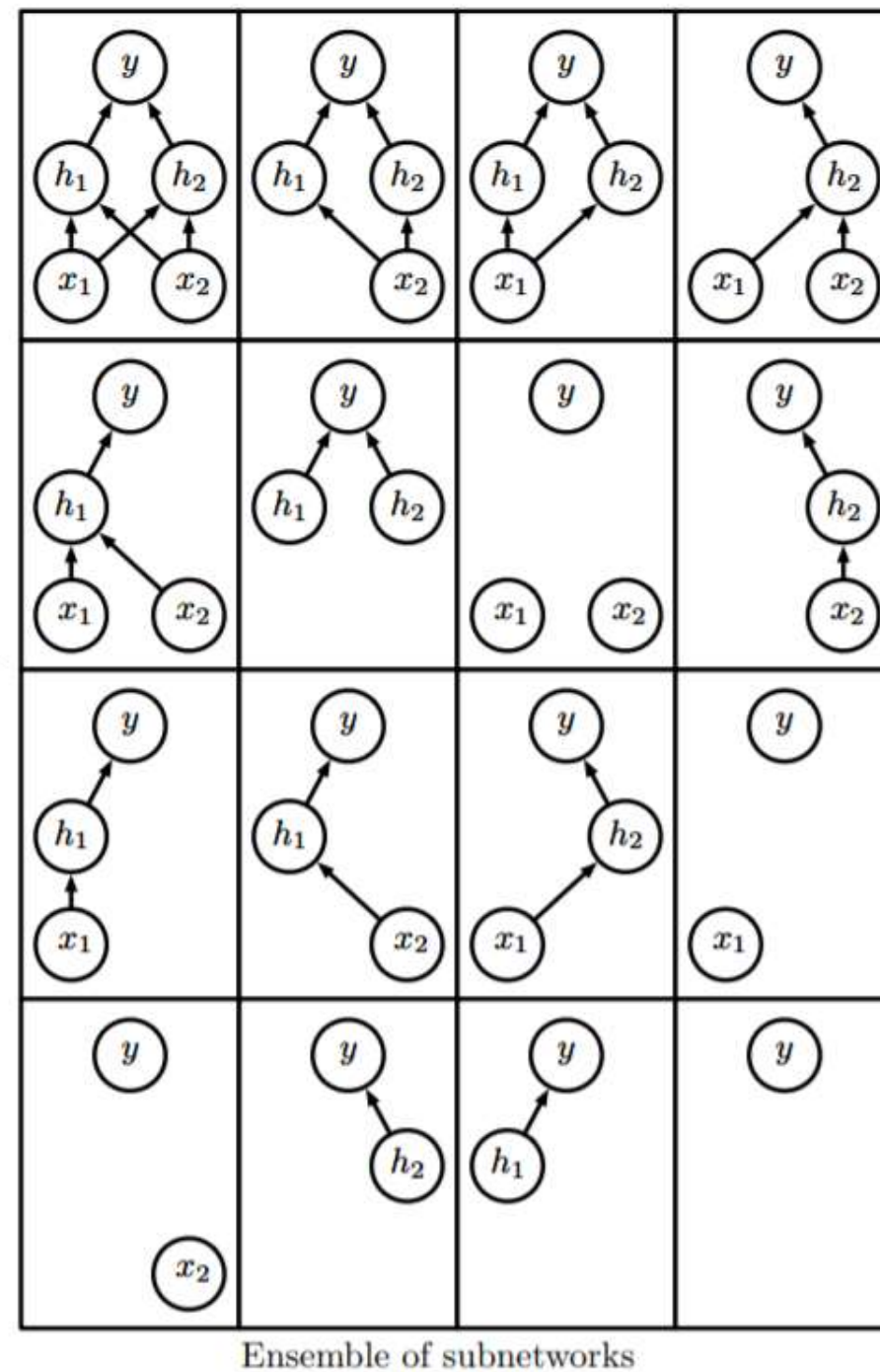
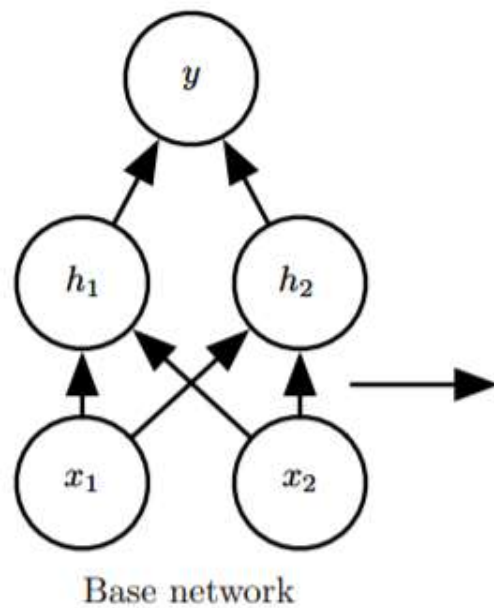
- Train  $k$  models.
- Suppose that the errors have variance  $v$  and covariance  $c$ .

$$\begin{aligned}\bar{e} &= \frac{1}{k} \sum_i e_i \\ \mathbb{E}[\bar{e}] &= \mathbb{E} \left[ \left( \frac{1}{k} \sum_i e_i \right)^2 \right] \\ &= \frac{1}{k^2} \mathbb{E} \left[ \left( \sum_i e_i \right)^2 \right] \\ &= \frac{1}{k^2} \mathbb{E} \left[ \sum_i \left( e_i^2 + \sum_{i \neq j} e_i e_j \right) \right] \\ &= \frac{1}{k^2} (k \mathbb{E}[e_i^2] + (k^2 - k) \mathbb{E}[e_i e_j]) \\ &= \frac{1}{k} \mathbb{E}[e_i^2] + \frac{k-1}{k} \mathbb{E}[e_i e_j] = \frac{1}{k} v + \frac{k-1}{k} c\end{aligned}$$

# Bagging/Model Ensemble



# Dropout



# Dropout: Implementation

The expected value of the neuron with dropout is:

$$E[ra] = \sum_r p(r)ra = pa + (1 - p)0 = pa$$

where  $a$  is o/p of neuron,  $r \in \{0,1\}$ , dropout probability  $p(r):p$

## Dropout

```
def train_step(X):
    hidden_layer_1 = np.maximum(0, np.dot(W1, X) + b1)
    dropout_mask_1 = np.random.binomial(1, keep_prob, hidden_layer_1.shape)
    hidden_layer_1 *= dropout_mask_1
    hidden_layer_2 = np.maximum(0, np.dot(W2, hidden_layer_1) + b2)
    dropout_mask_2 = np.random.binomial(1, keep_prob, hidden_layer_2.shape)
    hidden_layer_2 *= dropout_mask_2
    out = np.dot(W3, hidden_layer_2) + b3

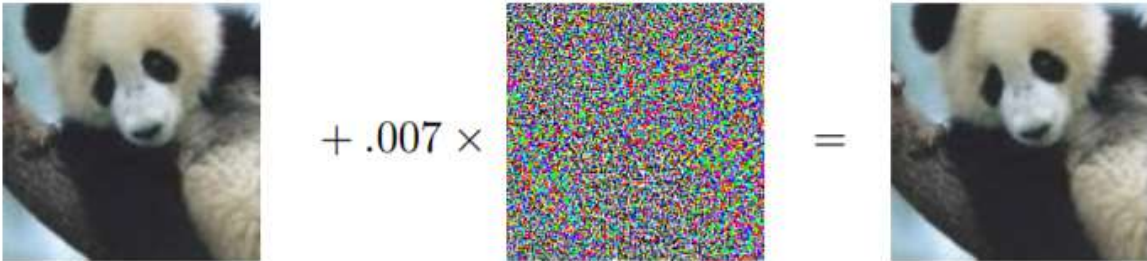
def predict(X):
    # ensembled forward pass
    hidden_layer_1 = np.maximum(0, np.dot(W1, X) + b1) * keep_prob
    hidden_layer_2 = np.maximum(0, np.dot(W2, hidden_layer_1) + b2) * keep_prob
    out = np.dot(W3, hidden_layer_2) + b3
```

## Inverted Dropout

```
def train_step(X):
    hidden_layer_1 = np.maximum(0, np.dot(W1, X) + b1)
    dropout_mask_1 = np.random.binomial(1, keep_prob, hidden_layer_1.shape) / keep_prob
    hidden_layer_1 *= dropout_mask_1
    hidden_layer_2 = np.maximum(0, np.dot(W2, hidden_layer_1) + b2)
    dropout_mask_2 = np.random.binomial(1, keep_prob, hidden_layer_2.shape) / keep_prob
    hidden_layer_2 *= dropout_mask_2
    out = np.dot(W3, hidden_layer_2) + b3

def predict(X):
    # ensembled forward pass
    hidden_layer_1 = np.maximum(0, np.dot(W1, X) + b1)
    hidden_layer_2 = np.maximum(0, np.dot(W2, hidden_layer_1) + b2)
    out = np.dot(W3, hidden_layer_2) + b3
```

# Adversarial Examples


$$\begin{array}{ccc} x & + .007 \times \text{sign}(\nabla_x J(\theta, x, y)) & = x + \epsilon \text{sign}(\nabla_x J(\theta, x, y)) \\ y = \text{"panda"} & \text{"nematode"} & \text{"gibbon"} \\ \text{w/ 57.7\% confidence} & \text{w/ 8.2\% confidence} & \text{w/ 99.3 \% confidence} \end{array}$$

Training on adversarial examples is mostly intended to improve security, but can sometimes provide generic regularization.



# Adversarial Examples Intuition

$$P(y = 1 \mid x; w, b) = \sigma(w^T x + b), \text{ where } \sigma(z) = 1/(1 + e^{-z})$$

Suppose

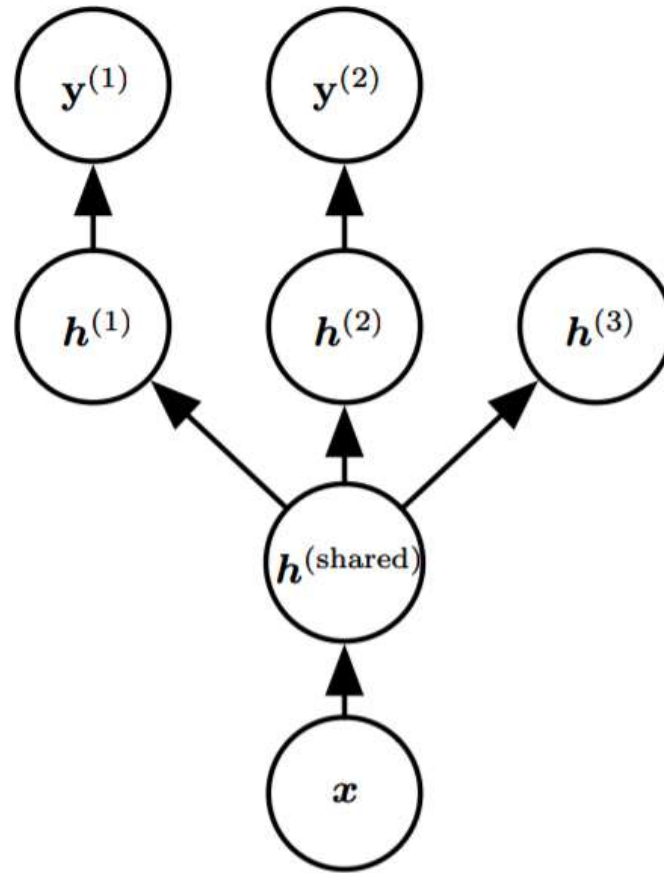
```
x = [2, -1, 3, -2, 2, 2, 1, -4, 5, 1] // input
w = [-1, -1, 1, -1, 1, -1, 1, 1, -1, 1] // weight vector
```

Hence, probability of class 1 is  $1/(1+e^{-( -3)}) = 0.0474$

We're now going to try to fool the classifier.

```
// xad = x + 0.5w gives:
xad = [1.5, -1.5, 3.5, -2.5, 2.5, 1.5, 1.5, -3.5, 4.5, 1.5]
```

# Multitask Learning



# Sparse representation:

## Direct versus Representational Sparsity

$$\begin{bmatrix} 18 \\ 5 \\ 15 \\ -9 \\ -3 \end{bmatrix} = \begin{bmatrix} 4 & 0 & 0 & -2 & 0 & 0 \\ 0 & 0 & -1 & 0 & 3 & 0 \\ 0 & 5 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & -1 & 0 & -4 \\ 1 & 0 & 0 & 0 & -5 & 0 \end{bmatrix} \begin{bmatrix} 2 \\ 3 \\ -2 \\ -5 \\ 1 \\ 4 \end{bmatrix} \quad (7.46)$$

$\mathbf{y} \in \mathbb{R}^m$                        $\mathbf{A} \in \mathbb{R}^{m \times n}$                        $\mathbf{x} \in \mathbb{R}^n$

$$\begin{bmatrix} -14 \\ 1 \\ 19 \\ 2 \\ 23 \end{bmatrix} = \begin{bmatrix} 3 & -1 & 2 & -5 & 4 & 1 \\ 4 & 2 & -3 & -1 & 1 & 3 \\ -1 & 5 & 4 & 2 & -3 & -2 \\ 3 & 1 & 2 & -3 & 0 & -3 \\ -5 & 4 & -2 & 2 & -5 & -1 \end{bmatrix} \begin{bmatrix} 0 \\ 2 \\ 0 \\ 0 \\ -3 \\ 0 \end{bmatrix} \quad (7.47)$$

$\mathbf{y} \in \mathbb{R}^m$                        $\mathbf{B} \in \mathbb{R}^{m \times n}$                        $\mathbf{h} \in \mathbb{R}^n$

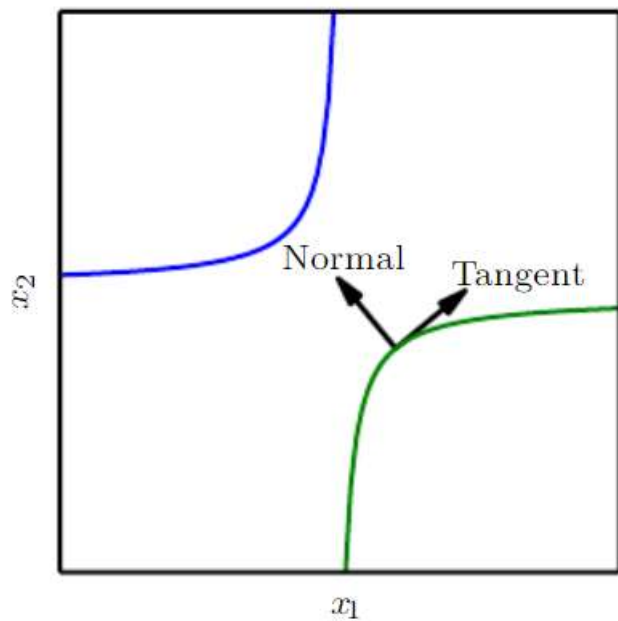
$$\tilde{J}(\boldsymbol{\theta}; \mathbf{X}, \mathbf{y}) = J(\boldsymbol{\theta}; \mathbf{X}, \mathbf{y}) + \alpha \Omega(\mathbf{h})$$

$$\Omega(\mathbf{h}) = \|\mathbf{h}\|_1 = \sum_i |h_i|$$

# Injecting noise

- Injecting noise into the input of a neural network can be seen as data augmentation
- Neural networks are not robust to noise
- To improve robustness, train them with random noise applied to their inputs
  - Part of some unsupervised learning, such as denoising autoencoder
- Noise can also be applied to hidden units
- Dropout, a powerful regularization strategy, can be viewed as constructing new inputs by multiplying by noise
- Noise applied to weights
- Injecting noise at the output targets

# Tangent prop algorithm



$$\Omega(f) = \sum_i \left( (\nabla_{\mathbf{x}} f(\mathbf{x}))^\top \mathbf{v}^{(i)} \right)^2$$



END