

# Regularization for Deep Learning

Lecture slides for Chapter 7 of *Deep Learning*

[www.deeplearningbook.org](http://www.deeplearningbook.org)

Ian Goodfellow

2016-09-27

Adapted by m.n. for CMPS 392

# Definition

- “Regularization is any modification we make to a learning algorithm that is intended to reduce its generalization error but not its training error.”
- Developing more effective regularization strategies has been one of the major research efforts in the field.
- Deep learning take:
  - the best fitting model (in the sense of minimizing generalization error) is a **large** model that has been **regularized** appropriately!

# Regularization strategies

- Constraints: adding restrictions on the parameter values.
- Soft constraints: Adding extra terms in the objective function:
  - Encode Prior knowledge.
  - Generic preference for a simpler model
- Ensemble methods:
  - Combine multiple hypotheses to explain the training data

# Parameter norm penalties

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

- $\theta$ : all learnable parameters (weights and biases)
- $w$ : parameters affected by a norm penalty
  - we take weights and exclude biases
- $\alpha \in [0, \infty)$ 
  - $\alpha = 0$ , no regularization
- $\Omega$ : norm function
  - L1
  - L2

# $L^2$ Parameter regularization aka. ridge regression

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

- $\nabla_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w} = \mathbf{w}$
- Update step:  $\mathbf{w} \leftarrow \mathbf{w} - \epsilon \alpha \mathbf{w} - \epsilon \nabla_{\mathbf{w}} J$ 
  - $\mathbf{w} \leftarrow \mathbf{w}(1 - \epsilon \alpha) - \epsilon \nabla_{\mathbf{w}} J$
- Let  $\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} J$
- Let  $\tilde{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}} \tilde{J}$
- Approximating  $J$  in the neighborhood of  $\mathbf{w}^*$ :
  - $\hat{J}(\mathbf{w}) = J(\mathbf{w}^*) + (\mathbf{w} - \mathbf{w}^*)^T \mathbf{H}(\mathbf{w} - \mathbf{w}^*)$

Weights are shrunk by a multiplicative factor

No first order term since  $\mathbf{w}^*$  is the minimum ( $\nabla J(\mathbf{w}^*) = \mathbf{0}$ )

# How $\tilde{\mathbf{w}}$ (regularized solution) compares to unregularized solution $\mathbf{w}^*$ ?

- What is the gradient of  $\hat{J}(\mathbf{w})$  at  $\tilde{\mathbf{w}}$ ?
- $\nabla \hat{J}(\tilde{\mathbf{w}}) = \mathbf{H}(\tilde{\mathbf{w}} - \mathbf{w}^*)$
- $\nabla \tilde{J}(\tilde{\mathbf{w}}) = \mathbf{H}(\tilde{\mathbf{w}} - \mathbf{w}^*) + \alpha\tilde{\mathbf{w}} = \mathbf{0}$ 
  - $(\mathbf{H} + \alpha\mathbf{I})\tilde{\mathbf{w}} = \mathbf{H}\mathbf{w}^*$
  - $\tilde{\mathbf{w}} = (\mathbf{H} + \alpha\mathbf{I})^{-1}\mathbf{H}\mathbf{w}^*$
- $\mathbf{H}$  is real and symmetric
  - $\mathbf{H} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T$
- $\tilde{\mathbf{w}} = (\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T + \alpha\mathbf{I})^{-1}(\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T)\mathbf{w}^* = (\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T + \mathbf{Q}\alpha\mathbf{I}\mathbf{Q}^T)^{-1}(\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T)\mathbf{w}^*$
- $\tilde{\mathbf{w}} = (\mathbf{Q}(\mathbf{\Lambda} + \alpha\mathbf{I})\mathbf{Q}^T)^{-1}(\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T)\mathbf{w}^* = \mathbf{Q}(\mathbf{\Lambda} + \alpha\mathbf{I})^{-1}\mathbf{Q}^T\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T\mathbf{w}^*$

$$\tilde{\mathbf{w}} = \mathbf{Q}(\mathbf{\Lambda} + \alpha\mathbf{I})^{-1}\mathbf{\Lambda}\mathbf{Q}^T\mathbf{w}^*$$

# Interpretation

$$\tilde{\mathbf{w}} = \mathbf{Q}(\mathbf{\Lambda} + \alpha \mathbf{I})^{-1} \mathbf{\Lambda} \mathbf{Q}^T \mathbf{w}^*$$

- $\mathbf{w}^*$  projections against the eigen vectors of  $\mathbf{H}$  are scaled
  - Component  $i$  is multiplied by  $\frac{\lambda_i}{\lambda_i + \alpha}$
  - $\lambda_i \gg \alpha \Rightarrow$  the effect of regularization is small
  - $\lambda_i \ll \alpha \Rightarrow$  the corresponding component is shrunk by a factor of  $\alpha$

# Weight Decay

$$\hat{J}(\mathbf{w}) = J(\mathbf{w}^*) + \frac{1}{2}(\mathbf{w} - \mathbf{w}^*)^T \mathbf{H}(\mathbf{w} - \mathbf{w}^*)$$

Small Eigen  
vector of  $\mathbf{H}$   
(regularization  
effect is large)

Unregularized  
solution

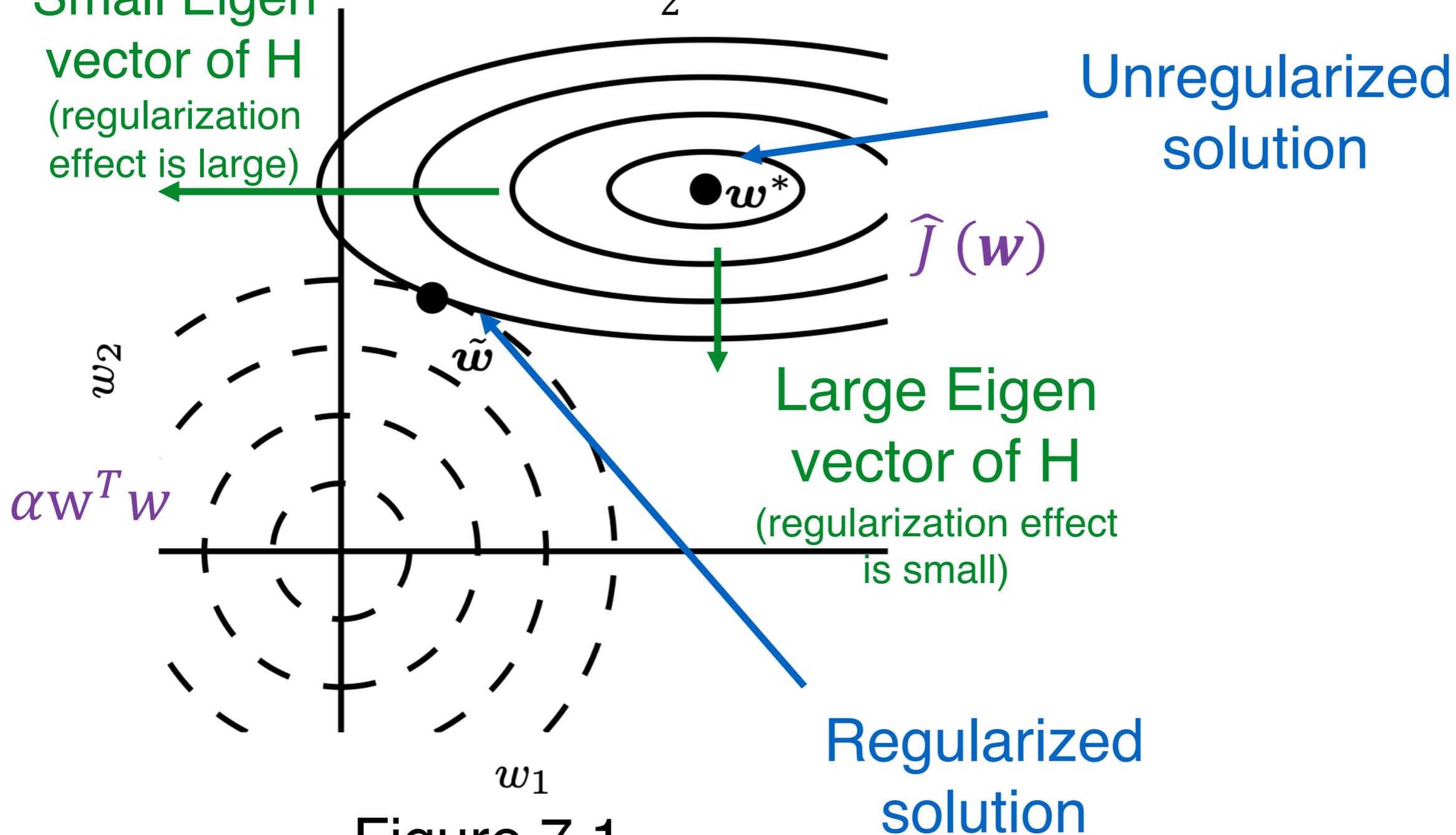


Figure 7.1

# Special case: Linear Regression

- Cost function:  $(X\mathbf{w} - \mathbf{y})^T (X\mathbf{w} - \mathbf{y}) + \frac{1}{2} \alpha \mathbf{w}^T \mathbf{w}$
- *Normal equations*
  - $X^T X \mathbf{w} - X^T \mathbf{y} + \alpha \mathbf{w} = \mathbf{0} \Rightarrow (X^T X + \alpha I) \mathbf{w} = X^T \mathbf{y}$
  - $\mathbf{w} = (X^T X + \alpha I)^{-1} X^T \mathbf{y}$ 
    - ← Covariance
    - ← feature-output
    - ← Proportional to the covariance matrix
- Basically, we are adding  $\alpha$  to the diag.
  - The diag. elements correspond to the variance of each feature
- We perceive the data as having higher variance
  - A feature having low covariance with output got shrunk even more due to this added variance

# $L^1$ regularization

- $\Omega(\boldsymbol{\theta}) = \|\mathbf{w}\|_1 = \sum_i |w_i|$
- $\tilde{J}(\mathbf{w}; \mathbf{X}, \mathbf{y}) = \alpha \|\mathbf{w}\|_1 + J(\mathbf{w}; \mathbf{X}, \mathbf{y})$ 
  - $\nabla_{\mathbf{w}} \tilde{J}(\mathbf{w}; \mathbf{X}, \mathbf{y}) = \alpha \text{sign}(\mathbf{w}) + \nabla_{\mathbf{w}} J(\mathbf{w}; \mathbf{X}, \mathbf{y})$
- $\hat{J}(\mathbf{w}) = J(\mathbf{w}^*) + \frac{1}{2} (\mathbf{w} - \mathbf{w}^*)^T \mathbf{H} (\mathbf{w} - \mathbf{w}^*)$ 
  - $\nabla \hat{J}(\mathbf{w}) = \mathbf{H} (\mathbf{w} - \mathbf{w}^*)$
- Assume that  $\mathbf{H} = \text{diag}([H_{1,1}, \dots, H_{n,n}])$ ,  $H_{i,i} > 0$ 
  - Linear regression after PCA
- $\tilde{J}(\mathbf{w}) \approx J(\mathbf{w}^*) + \sum_i \left[ \frac{1}{2} H_{i,i} (w_i - w_i^*)^2 + \alpha |w_i| \right]$
- Solution:  $w_i = \text{sign}(w_i^*) \max \left\{ |w_i^*| - \frac{\alpha}{H_{i,i}}, 0 \right\}$

# Interpretation

$$w_i = \text{sign}(w_i^*) \max \left\{ |w_i^*| - \frac{\alpha}{H_{i,i}}, 0 \right\}$$

- If  $w_i^* > 0$ :

- $w_i^* > \frac{\alpha}{H_{i,i}} \Rightarrow w_i$  is shifted towards 0 by  $\frac{\alpha}{H_{i,i}}$

- $w_i^* \leq \frac{\alpha}{H_{i,i}} \Rightarrow w_i = 0$

- If  $w_i^* < 0$ :

- $-w_i^* > \frac{\alpha}{H_{i,i}} \Rightarrow w_i = w_i + \frac{\alpha}{H_{i,i}}$

- $w_i$  is shifted towards 0 by  $\frac{\alpha}{H_{i,i}}$

- $-w_i^* \leq \frac{\alpha}{H_{i,i}} \Rightarrow w_i = 0$

# $L^1$ regularization sparsity

- The sparsity property induced by  $L^1$  regularization can be used as a feature selection mechanism
  - LASSO regression (least absolute shrinkage and selection operator)
- Equivalent to MAP Bayesian estimation with Laplace prior
  - the prior is an isotropic Laplace distribution over  $w \in \mathbb{R}^n$ :
    - Laplace  $\left(w_i; 0, \frac{1}{\alpha}\right) = \frac{1}{2\alpha} \exp(-\alpha |w_i|)$
    - $\log \text{Laplace} \left(w_i; 0, \frac{1}{\alpha}\right) = -\log 2\alpha - \alpha |w_i|$

$\log \text{posterior} \propto \log \text{likelihood} + \log \text{prior}$   
 $\max \log \text{posterior} \Leftrightarrow \min (\text{negative log likelihood} - \log \text{prior})$

# Norm Penalties

- MAP: Maximum A-Posteriori
- L1:
  - Encourages sparsity,
  - equivalent to MAP Bayesian estimation with Laplace prior
- Squared L2:
  - Encourages small weights,
  - equivalent to MAP Bayesian estimation with Gaussian prior

# Explicit constraints

- We want to constrain  $\Omega(\theta)$  to be less than some constant  $k$

□ construct a generalized Lagrange function

$$\mathcal{L}(\theta, \alpha; \mathbf{X}, \mathbf{y}) = J(\theta; \mathbf{X}, \mathbf{y}) + \alpha(\Omega(\theta) - k).$$

$$\theta^* = \arg \min_{\theta} \max_{\alpha, \alpha \geq 0} \mathcal{L}(\theta, \alpha).$$

- We can fix  $\alpha$  but lose  $k$

$$\theta^* = \arg \min_{\theta} \mathcal{L}(\theta, \alpha^*) = \arg \min_{\theta} J(\theta; \mathbf{X}, \mathbf{y}) + \alpha^* \Omega(\theta).$$

- The regularized training problem  $\tilde{J}$  is equivalent to the explicit constraints problem for an unknown  $k$ !

# Projection

- Sometimes we may wish to use explicit constraints rather than penalties.
  - we can modify algorithms such as stochastic gradient descent to take a step downhill on  $J(\theta)$  and then project  $\theta$  back to the nearest point that satisfies  $\Omega(\theta) < k$ .
- How to project?
  - Project into unit L2 ball:  $\mathbf{x} = \frac{\mathbf{y}}{\max\{1, \|\mathbf{y}\|_2\}}$ .
  - Project into unit L1 ball:
    - No closed-form solution
    - Numerical solution

# Dataset Augmentation

- Best way to regularize is to train with more data
  - create fake data and add it to the training set.
  - We can generate new  $(x, y)$  pairs easily just by transforming the  $x$  inputs in our training set.
  - particularly effective for *object recognition*
    - translating the training images a few pixels in each direction
    - rotating the image or scaling
- Some inappropriate transformations:
  - horizontal flips: 'b' and 'd',
  - 180° rotations: '6' and '9',

# Dataset Augmentation



Affine  
Distortion



Noise



Elastic  
Deformation



Horizontal  
flip



Random  
Translation



Hue Shift



# Noise Robustness

- Noise with infinitesimal variance can be added:
  - At the input
  - At the hidden layers
  - At the weights:  $\epsilon_{\mathbf{W}} \sim \mathcal{N}(\epsilon; \mathbf{0}, \eta \mathbf{I})$

$$\begin{aligned}\tilde{J}_{\mathbf{W}} &= \mathbb{E}_{p(\mathbf{x}, y, \epsilon_{\mathbf{W}})} \left[ (\hat{y}_{\epsilon_{\mathbf{W}}}(\mathbf{x}) - y)^2 \right] \\ &= \mathbb{E}_{p(\mathbf{x}, y, \epsilon_{\mathbf{W}})} \left[ \hat{y}_{\epsilon_{\mathbf{W}}}^2(\mathbf{x}) - 2y\hat{y}_{\epsilon_{\mathbf{W}}}(\mathbf{x}) + y^2 \right]\end{aligned}$$

# Injecting noise at the weights

$$\tilde{J}_{\mathbf{W}} = \mathbb{E}_{p(\mathbf{x}, y, \epsilon_{\mathbf{W}})} \left[ (\hat{y}_{\epsilon_{\mathbf{W}}}(\mathbf{x}) - y)^2 \right]$$

- For  $\eta$  small:
$$= \mathbb{E}_{p(\mathbf{x}, y, \epsilon_{\mathbf{W}})} \left[ \hat{y}_{\epsilon_{\mathbf{W}}}^2(\mathbf{x}) - 2y\hat{y}_{\epsilon_{\mathbf{W}}}(\mathbf{x}) + y^2 \right]$$
  - $\hat{y}_{\epsilon_{\mathbf{W}}} = \hat{y}(\mathbf{W} + \epsilon) = \hat{y}(\mathbf{W}) + \epsilon^T \nabla_{\mathbf{w}} \hat{y}(\mathbf{W})$
  - $\mathbb{E}_{p(\mathbf{x}, y, \epsilon_{\mathbf{W}})} [\hat{y}_{\epsilon_{\mathbf{W}}}^2] = \mathbb{E}_{p(\mathbf{x}, y)} [\hat{y}^2] + \mathbb{E}_{p(\epsilon_{\mathbf{W}})} [\epsilon^2] \mathbb{E}_{p(\mathbf{x}, y)} [\|\nabla_{\mathbf{w}} \hat{y}\|^2] + 0$
  - $\mathbb{E}_{p(\mathbf{x}, y, \epsilon_{\mathbf{W}})} [\hat{y}_{\epsilon_{\mathbf{W}}}^2] = \mathbb{E}_{p(\mathbf{x}, y)} [\hat{y}^2] + \eta \mathbb{E}_{p(\mathbf{x}, y)} [\|\nabla_{\mathbf{w}} \hat{y}\|^2]$
  - $\tilde{J}_{\mathbf{W}} = J + \eta \mathbb{E}_{p(\mathbf{x}, y)} [\|\nabla_{\mathbf{w}} \hat{y}(\mathbf{x})\|^2]$
  - Equivalent to adding a regularization term
- Pushes the model into regions where the model is relatively insensitive to small variations in the weights

# Special case: linear regression

- $\tilde{J}_{\mathbf{w}} = J + \eta \mathbb{E}_{p(\mathbf{x},y)} [\|\nabla_{\mathbf{w}} \hat{y}(\mathbf{x})\|^2]$
- $\hat{y} = \mathbf{w}^T \mathbf{x} + b$
- $\mathbb{E}_{p(\mathbf{x},y)} [\|\nabla_{\mathbf{w}} \hat{y}(\mathbf{x})\|^2] = \mathbb{E}_{p(\mathbf{x})} [\|\mathbf{x}\|^2]$
- which is not a function of parameters and therefore does not contribute to the cost function w.r.t  $\mathbf{w}$ :
  - No regularization effect!

# Injecting noise at the output targets

- Most datasets have some amount of mistakes in the  $y$  labels.
- It can be harmful to maximize  $\log p(y | x)$  when  $y$  is a mistake.
- One way to prevent this is to explicitly model the noise on the labels.
  - For example, we can assume that for some small constant  $\epsilon$ , the training set label  $y$  is correct with probability  $1 - \epsilon$ ,
  - and otherwise any of the other possible labels might be correct.
- This assumption is easy to incorporate into the cost function analytically,
  - rather than by explicitly drawing noise samples.
  - For example, **label smoothing** regularizes a model based on a softmax with  $k$  output values
    - by replacing the hard 0 by  $\frac{\epsilon}{k-1}$
    - and 1 by  $1 - \epsilon$
- Label smoothing has the advantage of preventing the pursuit of hard probabilities without discouraging correct classification.

# Multi-Task Learning

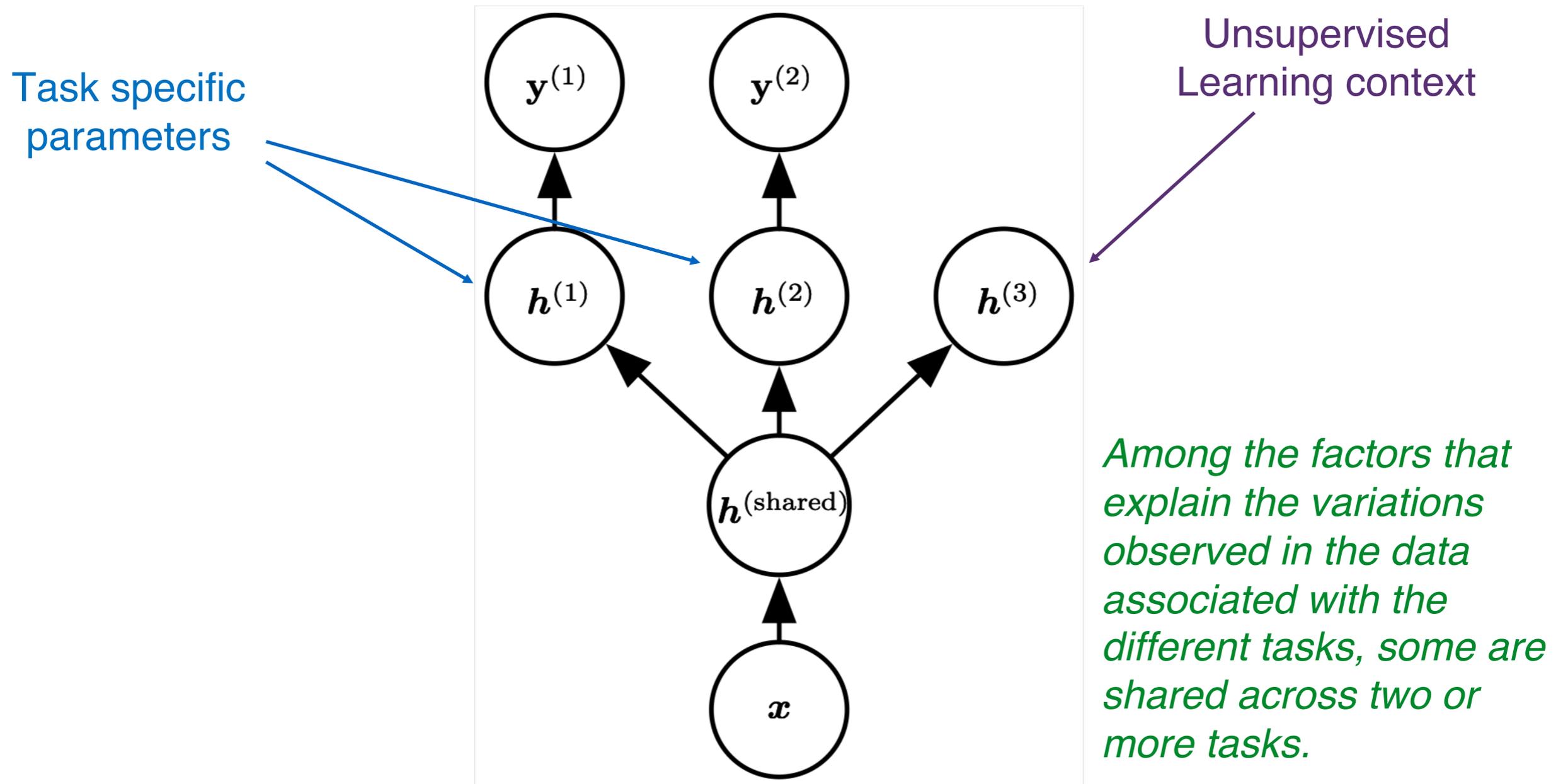


Figure 7.2

# Learning Curves

Early stopping: terminate while validation set performance is better

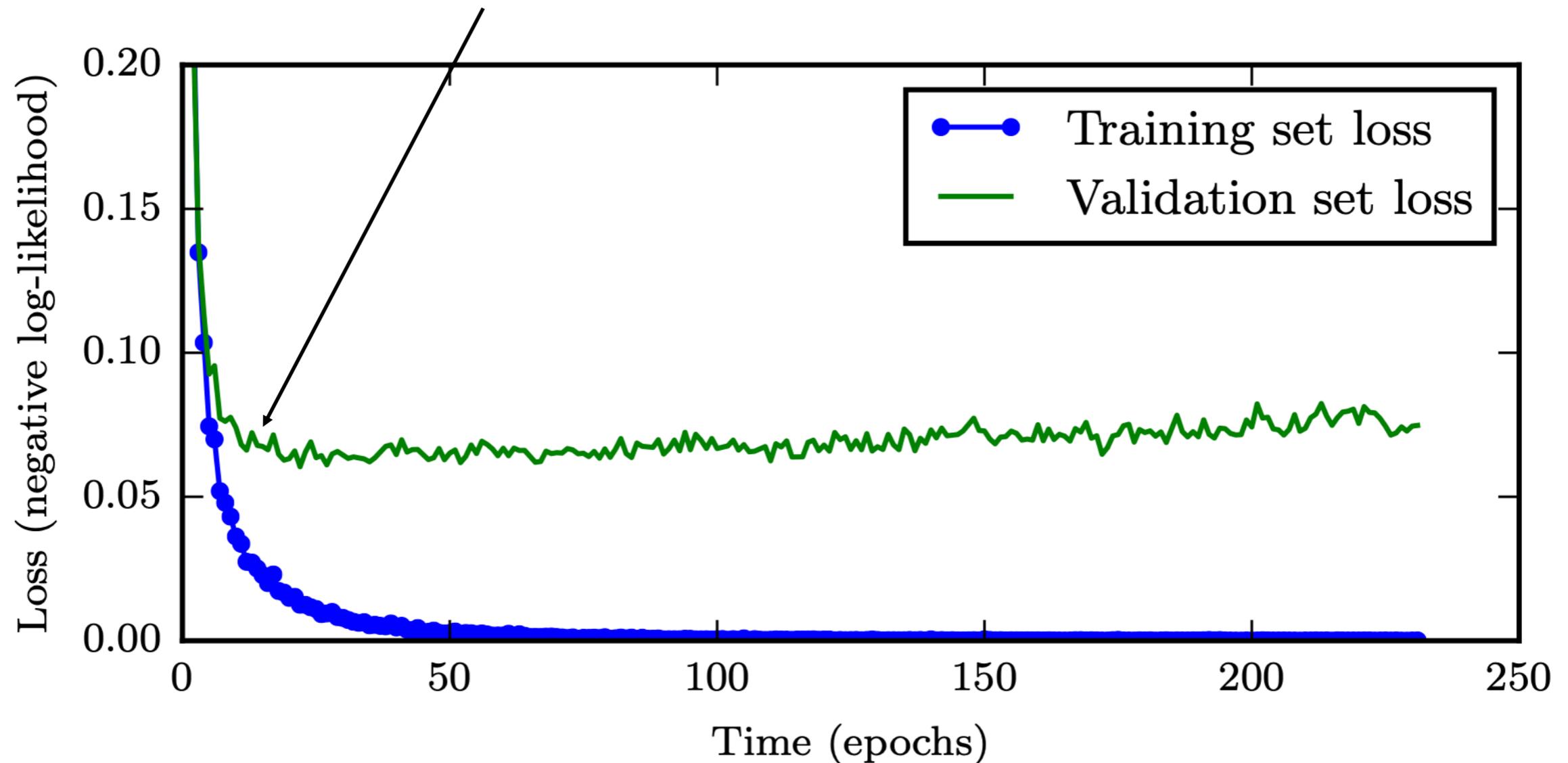


Figure 7.3

# Early stopping

- probably the most commonly used form of regularization in deep learning.
  - the number of training steps (or training time) is just another hyperparameter.
- The cost is running the validation set evaluation periodically during training
  - Reduce the validation set
  - Evaluate the validation loss less frequently
- Periodically save the trained model

# Early stopping algorithm

Let  $n$  be the number of steps between evaluations.

Let  $p$  be the “patience,” the number of times to observe worsening validation set error before giving up.

Let  $\theta_o$  be the initial parameters.

$\theta \leftarrow \theta_o$

$i \leftarrow 0$

$j \leftarrow 0$

$v \leftarrow \infty$

$\theta^* \leftarrow \theta$

$i^* \leftarrow i$

while  $j < p$  do

    Update  $\theta$  by running the training algorithm for  $n$  steps.

$i \leftarrow i + n$

$v' \leftarrow \text{ValidationSetError}(\theta)$

    if  $v' < v$  then

$j \leftarrow 0$

$\theta^* \leftarrow \theta$

$i^* \leftarrow i$

$v \leftarrow v'$

    else

$j \leftarrow j + 1$

    end if

end while

Best parameters are  $\theta^*$ , best number of training steps is  $i^*$

# Re-use the validation set

---

**Algorithm 7.2** A meta-algorithm for using early stopping to determine how long to train, then retraining on all the data.

---

Let  $\mathbf{X}^{(\text{train})}$  and  $\mathbf{y}^{(\text{train})}$  be the training set.

Split  $\mathbf{X}^{(\text{train})}$  and  $\mathbf{y}^{(\text{train})}$  into  $(\mathbf{X}^{(\text{subtrain})}, \mathbf{X}^{(\text{valid})})$  and  $(\mathbf{y}^{(\text{subtrain})}, \mathbf{y}^{(\text{valid})})$  respectively.

Run early stopping (algorithm 7.1) starting from random  $\theta$  using  $\mathbf{X}^{(\text{subtrain})}$  and  $\mathbf{y}^{(\text{subtrain})}$  for training data and  $\mathbf{X}^{(\text{valid})}$  and  $\mathbf{y}^{(\text{valid})}$  for validation data. This returns  $i^*$ , the optimal number of steps.

Set  $\theta$  to random values again.

Train on  $\mathbf{X}^{(\text{train})}$  and  $\mathbf{y}^{(\text{train})}$  for  $i^*$  steps.

---

---

**Algorithm 7.3** Meta-algorithm using early stopping to determine at what objective value we start to overfit, then continue training until that value is reached.

---

Let  $\mathbf{X}^{(\text{train})}$  and  $\mathbf{y}^{(\text{train})}$  be the training set.

Split  $\mathbf{X}^{(\text{train})}$  and  $\mathbf{y}^{(\text{train})}$  into  $(\mathbf{X}^{(\text{subtrain})}, \mathbf{X}^{(\text{valid})})$  and  $(\mathbf{y}^{(\text{subtrain})}, \mathbf{y}^{(\text{valid})})$  respectively.

Run early stopping (algorithm 7.1) starting from random  $\theta$  using  $\mathbf{X}^{(\text{subtrain})}$  and  $\mathbf{y}^{(\text{subtrain})}$  for training data and  $\mathbf{X}^{(\text{valid})}$  and  $\mathbf{y}^{(\text{valid})}$  for validation data. This updates  $\theta$ .

$\epsilon \leftarrow J(\theta, \mathbf{X}^{(\text{subtrain})}, \mathbf{y}^{(\text{subtrain})})$

while  $J(\theta, \mathbf{X}^{(\text{valid})}, \mathbf{y}^{(\text{valid})}) > \epsilon$  do

    Train on  $\mathbf{X}^{(\text{train})}$  and  $\mathbf{y}^{(\text{train})}$  for  $n$  steps.

end while

---

Less well-  
behaved  $\longrightarrow$

# Early stopping as a regularizer

- $\epsilon$  (learning rate) and  $\tau$  (number of training steps) limits the the volume of parameters reachable from  $\theta_0$  (initial parameters)
- Early stopping is equivalent to L2 regularization in the case of:
  - a simple linear model
  - with a quadratic error function
  - and simple gradient descent
- $\hat{J}(\mathbf{w}) = J(\mathbf{w}^*) + \frac{1}{2}(\mathbf{w} - \mathbf{w}^*)^T \mathbf{H}(\mathbf{w} - \mathbf{w}^*)$ 
  - $\nabla_{\mathbf{w}} \hat{J}(\mathbf{w}) = \mathbf{H}(\mathbf{w} - \mathbf{w}^*)$
- $\mathbf{w}^{(\tau)} = \mathbf{w}^{(\tau-1)} - \epsilon \nabla_{\mathbf{w}} \hat{J}(\mathbf{w}^{(\tau-1)}) = \mathbf{w}^{(\tau-1)} - \epsilon \mathbf{H}(\mathbf{w}^{(\tau-1)} - \mathbf{w}^*)$ 
  - $\mathbf{w}^{(\tau)} = (\mathbf{I} - \epsilon \mathbf{H})\mathbf{w}^{(\tau-1)} + \epsilon \mathbf{H}\mathbf{w}^*$
  - $\mathbf{w}^{(\tau)} - \mathbf{w}^* = (\mathbf{I} - \epsilon \mathbf{H})\mathbf{w}^{(\tau-1)} + (\epsilon \mathbf{H} - \mathbf{I})\mathbf{w}^*$
  - $\mathbf{w}^{(\tau)} - \mathbf{w}^* = (\mathbf{I} - \epsilon \mathbf{H})(\mathbf{w}^{(\tau-1)} - \mathbf{w}^*)$

The number of steps  $\tau$  corresponds to some value of the weight decay coefficient  $\alpha$

- $\mathbf{w}^{(\tau)} - \mathbf{w}^* = (\mathbf{I} - \epsilon \mathbf{H})(\mathbf{w}^{(\tau-1)} - \mathbf{w}^*)$ 
  - $\mathbf{H} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T$
  - $\mathbf{w}^{(\tau)} - \mathbf{w}^* = \mathbf{Q}(\mathbf{I} - \epsilon \mathbf{\Lambda})\mathbf{Q}^T(\mathbf{w}^{(\tau-1)} - \mathbf{w}^*)$
  - $\mathbf{Q}^T(\mathbf{w}^{(\tau)} - \mathbf{w}^*) = (\mathbf{I} - \epsilon \mathbf{\Lambda})\mathbf{Q}^T(\mathbf{w}^{(\tau-1)} - \mathbf{w}^*)$ 
    - if  $\epsilon$  small  $\Rightarrow |1 - \epsilon \lambda_i| < 1$ , every step brings closer to  $\mathbf{w}^*$
- Assume we start with  $\mathbf{w}^{(0)} = \mathbf{0}$ :
  - $\mathbf{Q}^T \mathbf{w}^{(1)} = (\mathbf{I} - (\mathbf{I} - \epsilon \mathbf{\Lambda}))\mathbf{Q}^T \mathbf{w}^*$
  - $\mathbf{Q}^T \mathbf{w}^{(2)} = (\mathbf{I} - (\mathbf{I} - \epsilon \mathbf{\Lambda})^2)\mathbf{Q}^T \mathbf{w}^*$
  - $\mathbf{Q}^T \mathbf{w}^{(\tau)} = (\mathbf{I} - (\mathbf{I} - \epsilon \mathbf{\Lambda})^\tau)\mathbf{Q}^T \mathbf{w}^*$
- L2 regularization:
  - $\mathbf{Q}^T \tilde{\mathbf{w}} = (\mathbf{\Lambda} + \alpha \mathbf{I})^{-1} \mathbf{\Lambda} \mathbf{Q}^T \mathbf{w}^*$
  - $\mathbf{Q}^T \tilde{\mathbf{w}} = \mathbf{I} - (\mathbf{\Lambda} + \alpha \mathbf{I})^{-1} \alpha \mathbf{Q}^T \mathbf{w}^*$ 

$\frac{\lambda_i}{\lambda_i + \alpha} = 1 - \frac{\alpha}{\lambda_i + \alpha}$
- Compare:
  - $\mathbf{I} - (\mathbf{\Lambda} + \alpha \mathbf{I})^{-1} \alpha$  and  $\mathbf{I} - (\mathbf{I} - \epsilon \mathbf{\Lambda})^\tau$ 

$1 - (1 - \epsilon \lambda_i)^\tau$
  - $(\mathbf{\Lambda} + \alpha \mathbf{I})^{-1} \alpha = (\mathbf{I} - \epsilon \mathbf{\Lambda})^\tau$

# Early stopping advantage

- $\frac{\alpha}{\lambda_i + \alpha} = (1 - \epsilon\lambda_i)^\tau \Rightarrow \tau \log(1 - \epsilon\lambda_i) = \log\left(\frac{1}{1 + \frac{\lambda_i}{\alpha}}\right)$
- Assume  $\log(1 + x) \approx x$  for small enough  $x$ 
  - Assume  $\frac{\lambda_i}{\alpha} \ll 1$  and  $\epsilon\lambda_i \ll 1$
- $-\tau\epsilon\lambda_i \approx -\frac{\lambda_i}{\alpha} \Rightarrow \alpha \approx \frac{1}{\tau\epsilon}$ 
  - the number of training iterations  $\tau$  plays a role inversely proportional to the L2 regularization parameter,
  - and the inverse of  $\tau\epsilon$  plays the role of the weight decay coefficient.
- Early stopping advantage over weight decay:
  - early stopping automatically determines the correct amount of regularization
  - while weight decay requires many training experiments with different values of its hyperparameter.

# Early Stopping and Weight Decay

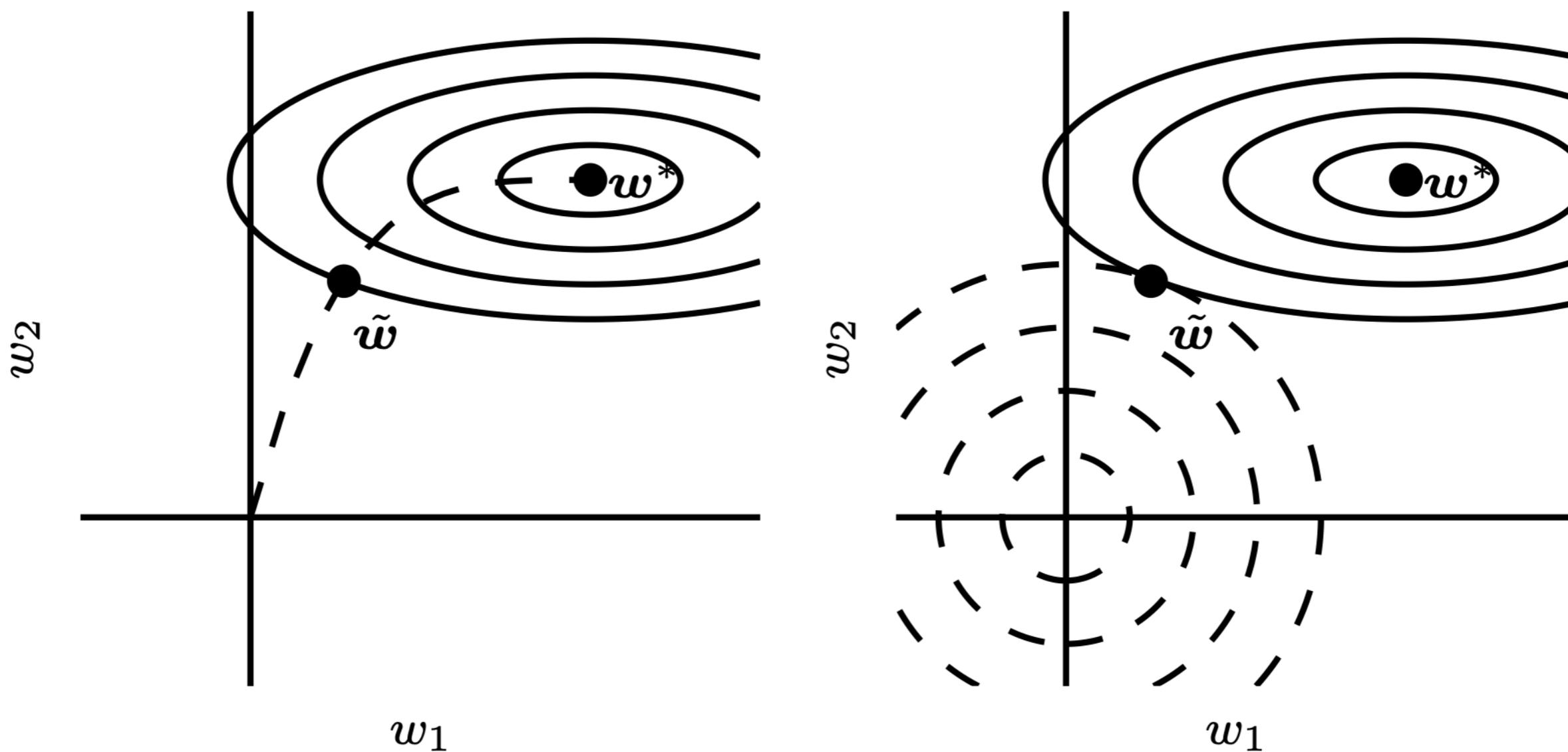


Figure 7.4

# Parameter tying

- Formally, we have model  $A$  with parameters  $\mathbf{w}^{(A)}$  and model  $B$  with parameters  $\mathbf{w}^{(B)}$
- The two models map the input to two different, but related outputs:
  - $\hat{y}^{(A)} = f(\mathbf{w}^{(A)}, \mathbf{x})$
  - $\hat{y}^{(B)} = g(\mathbf{w}^{(B)}, \mathbf{x})$
  - $\forall i, \mathbf{w}^{(A)}$  should be close to  $\mathbf{w}^{(B)}$
- Regularization

$$\Omega(\mathbf{w}^{(A)}, \mathbf{w}^{(B)}) = \|\mathbf{w}^{(A)} - \mathbf{w}^{(B)}\|_2^2$$

# Parameter sharing (e.g. CNN)

- Force sets of parameters to be equal.
- Advantage:
  - only a subset of the parameters (the unique set) needs to be stored in memory.
- Natural images have many statistical properties that are invariant to translation.
  - a photo of a cat remains a photo of a cat if it is translated one pixel to the right
  - Parameter sharing has allowed CNNs to dramatically lower the number of unique model parameters

# Sparse Representations

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

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

Sparse  
parameters

$$\Omega(\boldsymbol{\theta}) = \|\mathbf{w}\|_1 = \sum_i |w_i|$$

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

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

Sparse  
representations

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

# Bagging

- Bagging (short for bootstrap aggregating) is a technique for reducing generalization error by combining several models
  - train several different models separately
  - the models vote on the output for test examples
- Bagging is an example of model averaging.
  - The general term is Ensemble methods.
- The reason that model averaging works is that different models will usually not make all the same errors on the test set.

# Bagging example

- Consider for example a set of  $k$  regression models.
- Suppose that each model makes an error  $\epsilon_i$  on each example, with the errors drawn from a zero-mean multivariate normal distribution
  - with variances  $E[\epsilon_i^2] = v$
  - and covariances  $E[\epsilon_i \epsilon_j] = c$
- Then the error made by the average prediction of all the ensemble models is  $\frac{1}{k} \sum_i \epsilon_i$
- The expected squared error of the ensemble predictor is:

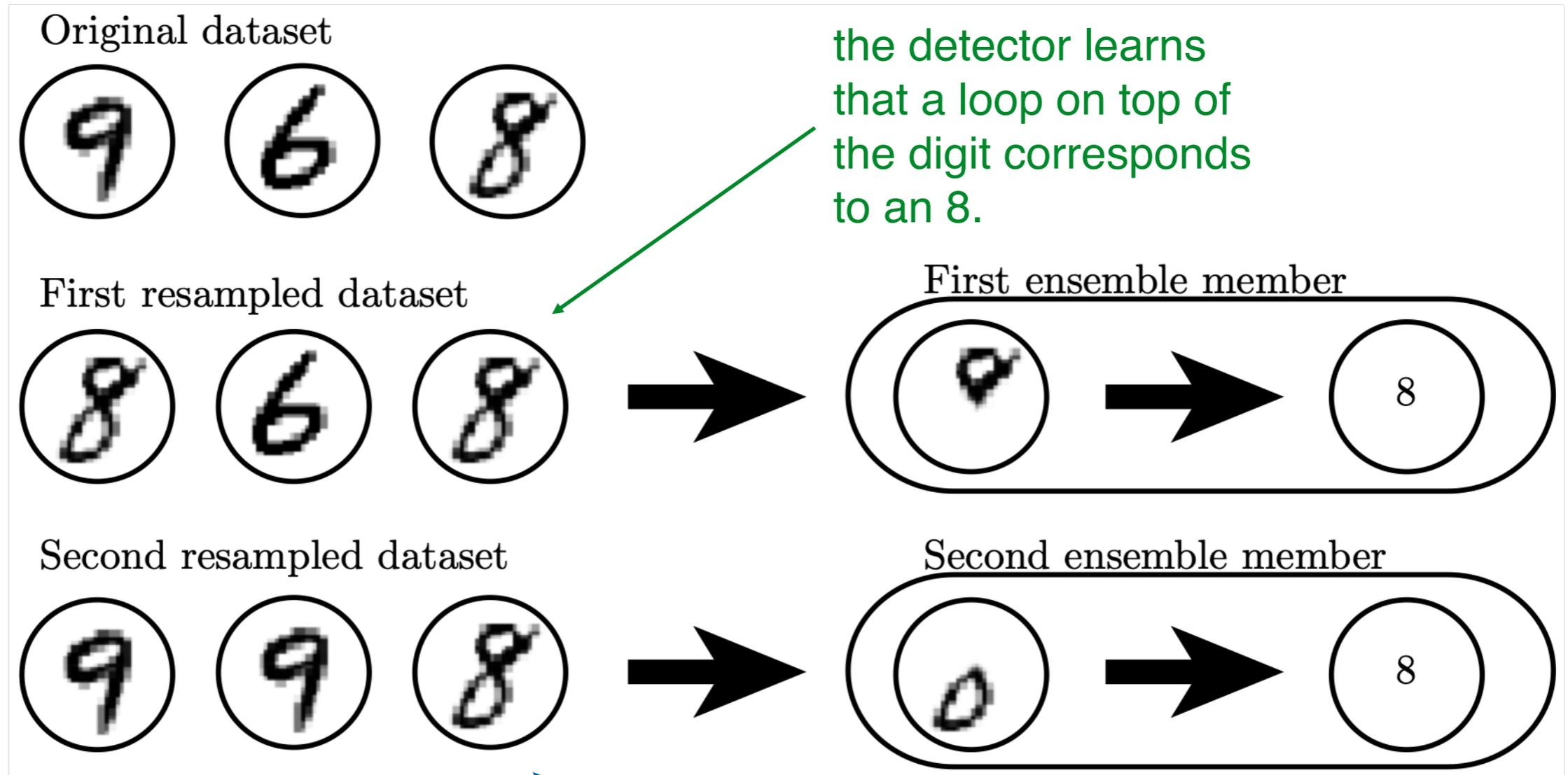
$$\begin{aligned} \mathbb{E} \left[ \left( \frac{1}{k} \sum_i \epsilon_i \right)^2 \right] &= \frac{1}{k^2} \mathbb{E} \left[ \sum_i \left( \epsilon_i^2 + \sum_{j \neq i} \epsilon_i \epsilon_j \right) \right] \\ &= \frac{1}{k} v + \frac{k-1}{k} c. \end{aligned}$$

- $c = v \Rightarrow$  no gain, the expected error remains  $v$
- $c = 0 \Rightarrow$  max gain, the expected error is  $v/k$

# Ensemble methods vs. bagging

- Different ensemble methods construct the ensemble of models in different ways.
- Bagging is a method that allows the same kind of model, training algorithm and objective function to be reused several times
- Bagging involves constructing  $k$  different datasets.
  - Each dataset has the same number of examples as the original dataset,
  - but each dataset is constructed by sampling with replacement from the original dataset.
    - with high probability, each dataset is missing some of the examples from the original dataset and also contains several duplicate examples
    - on average around  $2/3$  of the examples from the original dataset are found in the resulting training set, if it has the same size as the original

# Bagging



# Why 2/3 ?

- *N*: number of items
- *K*: number of unique items
- *A*: number of drawn items

$$P(k) = \frac{N!}{(N-k)! N^A} \left\{ \begin{matrix} A \\ k \end{matrix} \right\}$$

All permutations of *k* among *N* items

All ways to distribute *A* items among *k* subsets such as no subset is left empty

All possible ways to draw *A* items

- $\left\{ \begin{matrix} A \\ k \end{matrix} \right\}$  is a Stirling number of the second kind

# Expected number of duplicates

- The indicator  $d_i$  corresponds to original item  $i$ , taking the value of one if  $i$  is present and zero if not
- $P(d_i = 0) = 1 - \left(\frac{1}{N}\right)^A$
- $E[d_i] = 1 - \left(1 - \frac{1}{N}\right)^A$
- $E[\sum d_i] = \sum E[d_i] = NE[d_i] = N \left(1 - \left(1 - \frac{1}{N}\right)^A\right)$
- $A = N \Rightarrow E[k] = N \left(1 - \left(1 - \frac{1}{N}\right)^N\right) \rightarrow N(1 - e^{-1})$
- $E[k] \approx 0.632 N$

# More about bagging

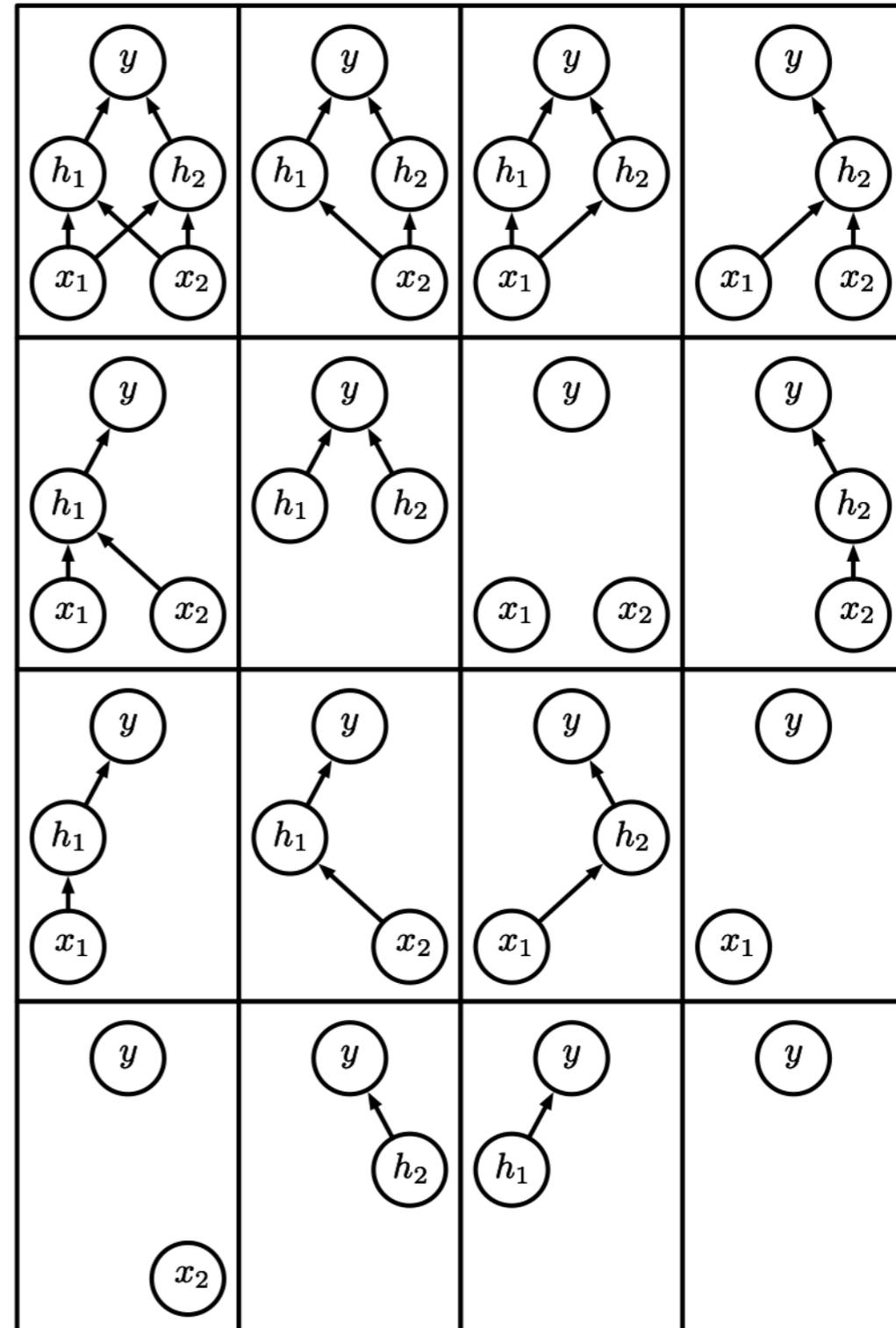
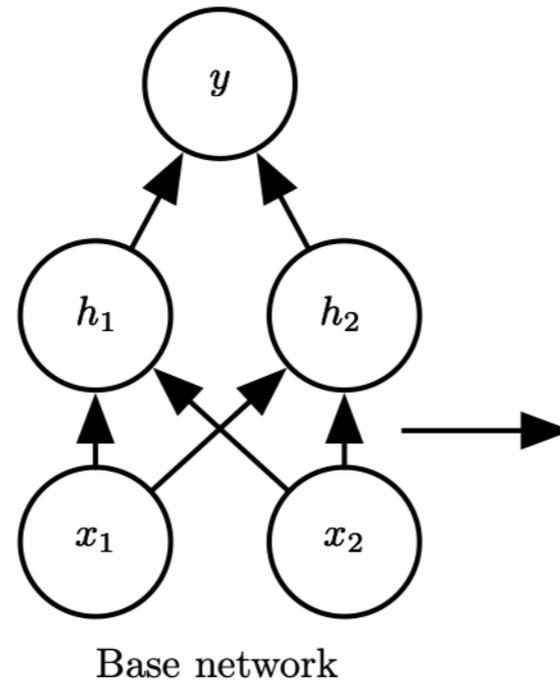
- Neural networks reach a wide enough variety of solution points that they can often benefit from model averaging
- Model averaging is an extremely powerful and reliable method for reducing generalization error.
  - Its use is usually discouraged when benchmarking algorithms for scientific papers
- Machine learning contests are usually won by methods using model averaging over dozens of models.

# Dropout

- Dropout provides an inexpensive approximation to training and evaluating a bagged ensemble of exponentially many neural networks.
  - removing non-output units from an underlying base network
    - by multiplying its output value by zero
- Each time we load an example into a minibatch, we randomly sample a different binary mask to apply to all of the input and hidden units in the network.
  - The probability of sampling a mask value of one (causing a unit to be included) is a hyperparameter fixed before training begins.
    - Typically, an input unit is included with probability 0.8 and a hidden unit is included with probability 0.5

# Dropout

Figure 7.6

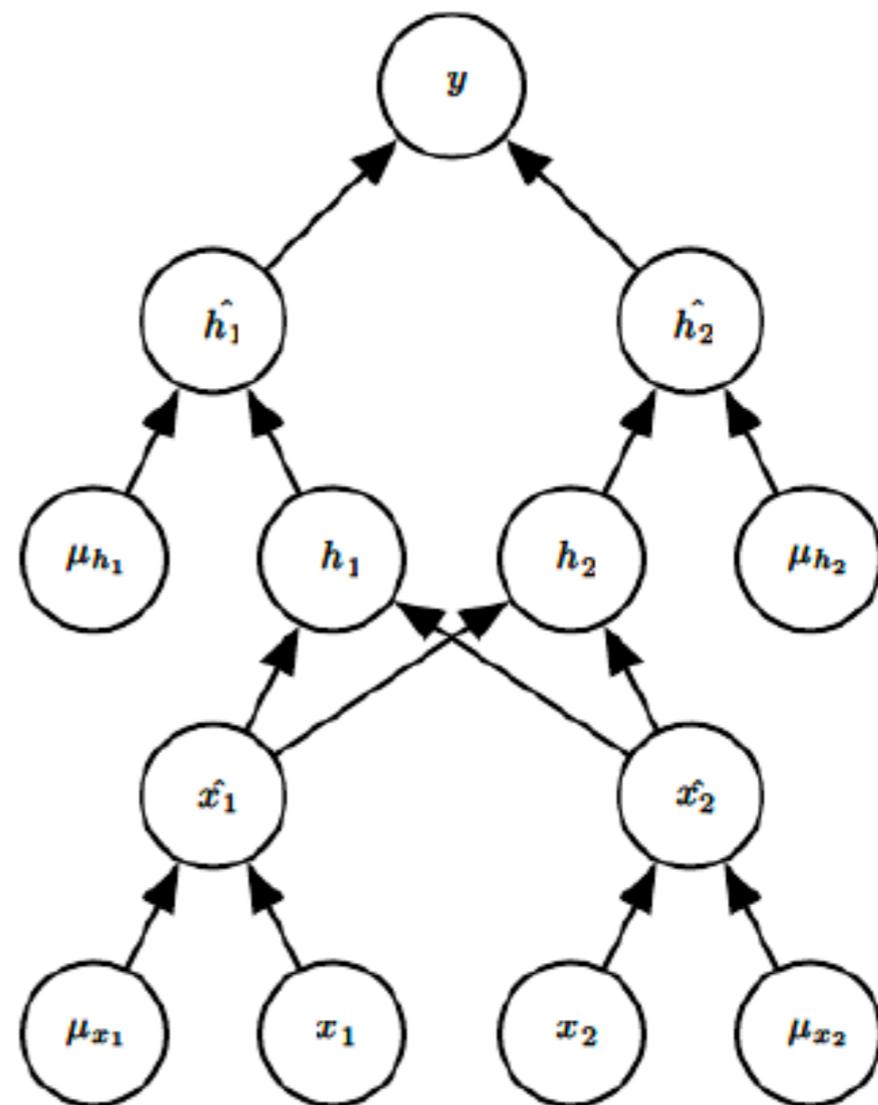
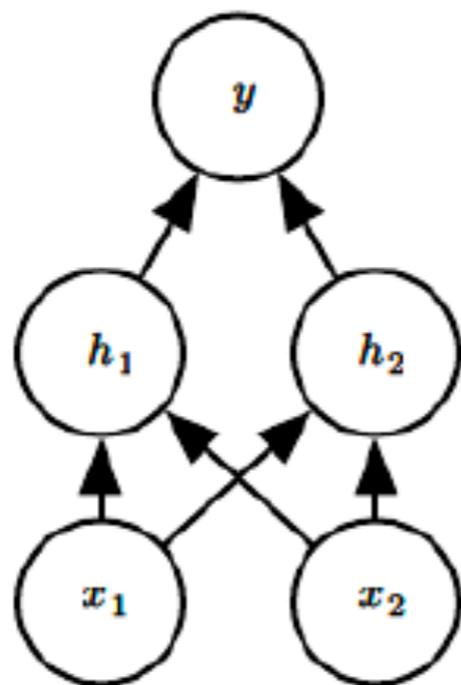


In networks with wider layers, the probability of dropping all possible paths from inputs to outputs becomes smaller.

# Dropout vs. bagging

- More formally, suppose that a mask vector  $\mu$  specifies which units to include, and  $J(\theta, \mu)$  defines the cost of the model defined by parameters  $\theta$  and mask  $\mu$ .
  - Then dropout training consists in minimizing  $\mathbb{E}_{\mu} J(\theta, \mu)$ .
  - The expectation contains exponentially many terms ( $2^d$ )
- Dropout training is not quite the same as bagging training.
  - In the case of bagging, the models are all independent.
  - In the case of dropout, the models share parameters
  - In bagging, each model is trained to convergence on its respective training set
  - In dropout, a tiny fraction of the possible sub-networks are each trained for a single step
  - In both, the training set encountered by each sub-network is a subset of the original training set sampled with replacement

# Computational graph of dropout



- The entries of  $\mu$  are binary and are sampled independently from each other,
  - And is not a function of the current value of the model parameters or the input example

# Inference

- To make a prediction, a bagged ensemble must accumulate votes from all of its members.
  - We refer to this process as inference
- In bagging, the prediction of the ensemble is  $\frac{1}{k} \sum_{i=1}^k p^{(i)}(y|\mathbf{x})$
- In dropout, the arithmetic mean is  $\sum_{\mu} p(\mu)p(y|\mathbf{x}, \mu)$
- The geometric mean is

$$\tilde{p}_{\text{ensemble}}(y|\mathbf{x}) = \sqrt[2^n]{\prod_{\mu} p(y|\mathbf{x}, \mu)}$$

- To guarantee that the result is a probability distribution,
  - we impose that none of the sub-models assigns probability 0 to any event,
  - and we renormalize the resulting distribution.

# Weight scaling inference rule

- Evaluate with the trained model with all units,
  - But with the weights going out of unit  $i$  multiplied by the probability of including unit  $i$  (e.g.  $\frac{1}{2}$ )
  - This corresponds to predict the geometric mean of the ensemble!
- Consider a softmax regression classifier with  $n$  input variables represented by the vector  $\mathbf{v}$ :

$$P(y = y_i | \mathbf{v}) = \text{softmax}(\mathbf{W}^T \mathbf{v} + \mathbf{b})_i$$

- To index into the family of submodels:

$$P(y = y_i | \mathbf{v}) = \text{softmax}(\mathbf{W}^T (\mathbf{v} \odot \mathbf{d}) + \mathbf{b})_i$$

$$\tilde{p}_{\text{ensemble}}(y = y_i | \mathbf{x}) = \sqrt[2^n]{\prod_{\mathbf{d} \in \{0,1\}^n} \text{softmax}(\mathbf{W}^T (\mathbf{v} \odot \mathbf{d}) + \mathbf{b})_i}$$

$$\tilde{p}_{\text{ensemble}}(y = y_i | \mathbf{x}) \propto \exp\left(\frac{1}{2^n} \sum_{\mathbf{d} \in \{0,1\}^n} (\mathbf{W}^T (\mathbf{v} \odot \mathbf{d}) + \mathbf{b})_i\right)$$

$$= \exp\left(\frac{1}{2^n} (2^{n-1} \mathbf{W}^T \mathbf{v} + 2^n \mathbf{b})_i\right) = \exp\left(\frac{1}{2} \mathbf{W}^T \mathbf{v} + \mathbf{b}\right)_i$$

# Another perspective of dropout

- (1) Dropout is bagging with parameter sharing
- (2) Information erasing: Each hidden unit must be able to perform well regardless of which other hidden units are in the model
  - ❑ Dropout thus regularizes each hidden unit to be not merely a good feature but a feature that is good in many contexts.
  - ❑ For example, if the model learns a hidden unit  $h_i$  that detects a face by finding the nose,
  - ❑ then dropping  $h_i$  corresponds to erasing the information that there is a nose in the image.
  - ❑ The model must learn another  $h_i$ ,
    - either that redundantly encodes the presence of a nose,
    - or that detects the face by another feature, such as the mouth

# Adversarial examples

- Search for an input  $x'$  near a data point  $x$  such that the model output is very different at  $x'$
- In many cases,  $x'$  can be so similar to  $x$  that a human observer cannot tell the difference between the original example and the adversarial example,
  - but the network can make highly different predictions.
- **Adversarial training**
  - training on adversarially perturbed examples from the training set
- Adversarial examples are interesting in the context of **regularization**
  - because one can reduce the error rate on the original i.i.d. test set via **adversarial training**

# Adversarial Examples

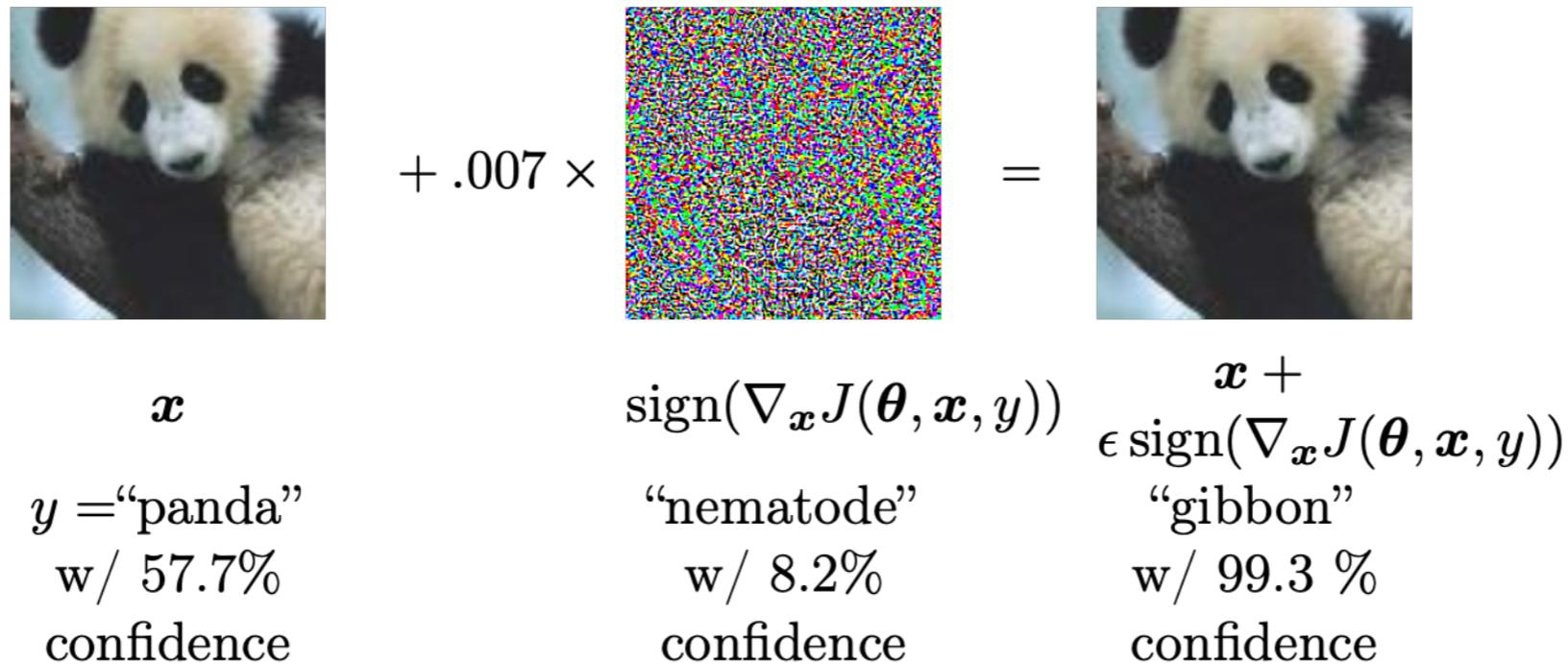


Figure 7.8

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

# Aversarial training

- The value of a linear function can change very rapidly if it has numerous inputs.
  - If we change each input by  $\epsilon$ , then a linear function with weights  $w$  can change by as much as  $\epsilon\|w\|_1$ , which can be a very large amount if  $w$  is high-dimensional.
- Adversarial training discourages this highly sensitive locally linear behavior by encouraging the network to be locally constant in the neighborhood of the training data.
- This can be seen as a way of explicitly introducing a local constancy prior into supervised neural nets.
  - The classifier may then be trained to assign the same label to  $x$  and  $x'$ .
  - The assumption motivating this approach is that different classes usually lie on disconnected manifolds, and a small perturbation should not be able to jump from one class manifold to another class manifold.

# Conclusion

- This chapter has described most of the general strategies used to regularize neural networks.
- Regularization is a central theme of machine learning

Our next topic is: **optimization**