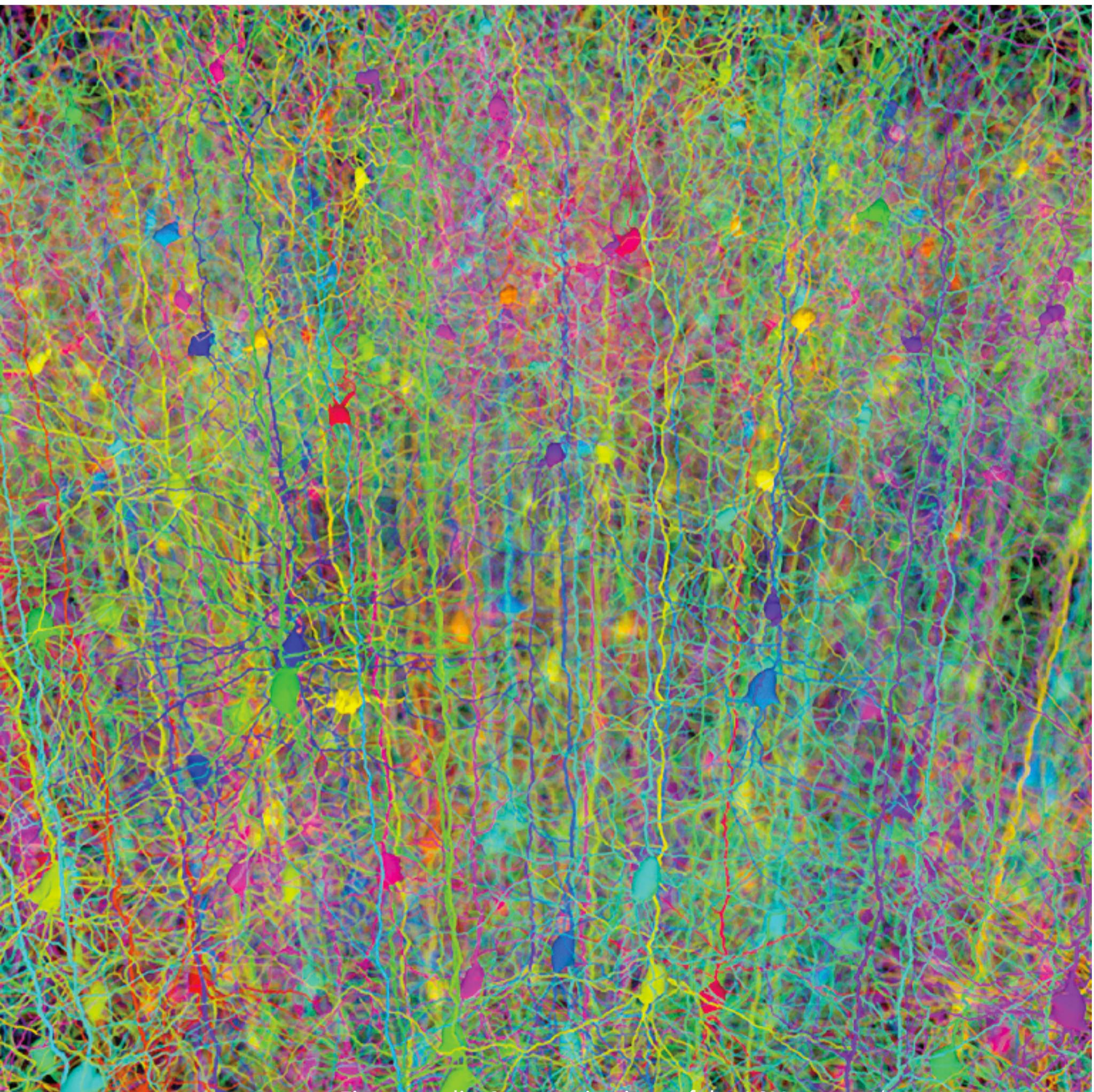


Regularization

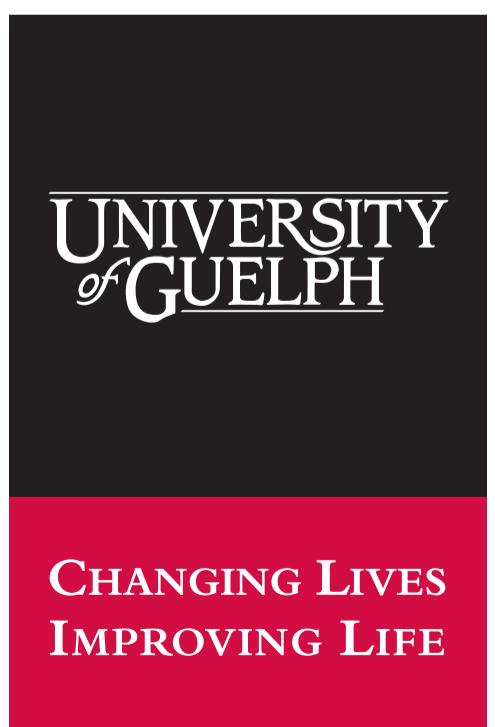


GRAHAM TAYLOR

VECTOR INSTITUTE

SCHOOL OF ENGINEERING
UNIVERSITY OF GUELPH

CANADIAN INSTITUTE
FOR ADVANCED RESEARCH



CIFAR
CANADIAN
INSTITUTE
FOR
ADVANCED
RESEARCH

Bias-Variance Trade-off

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



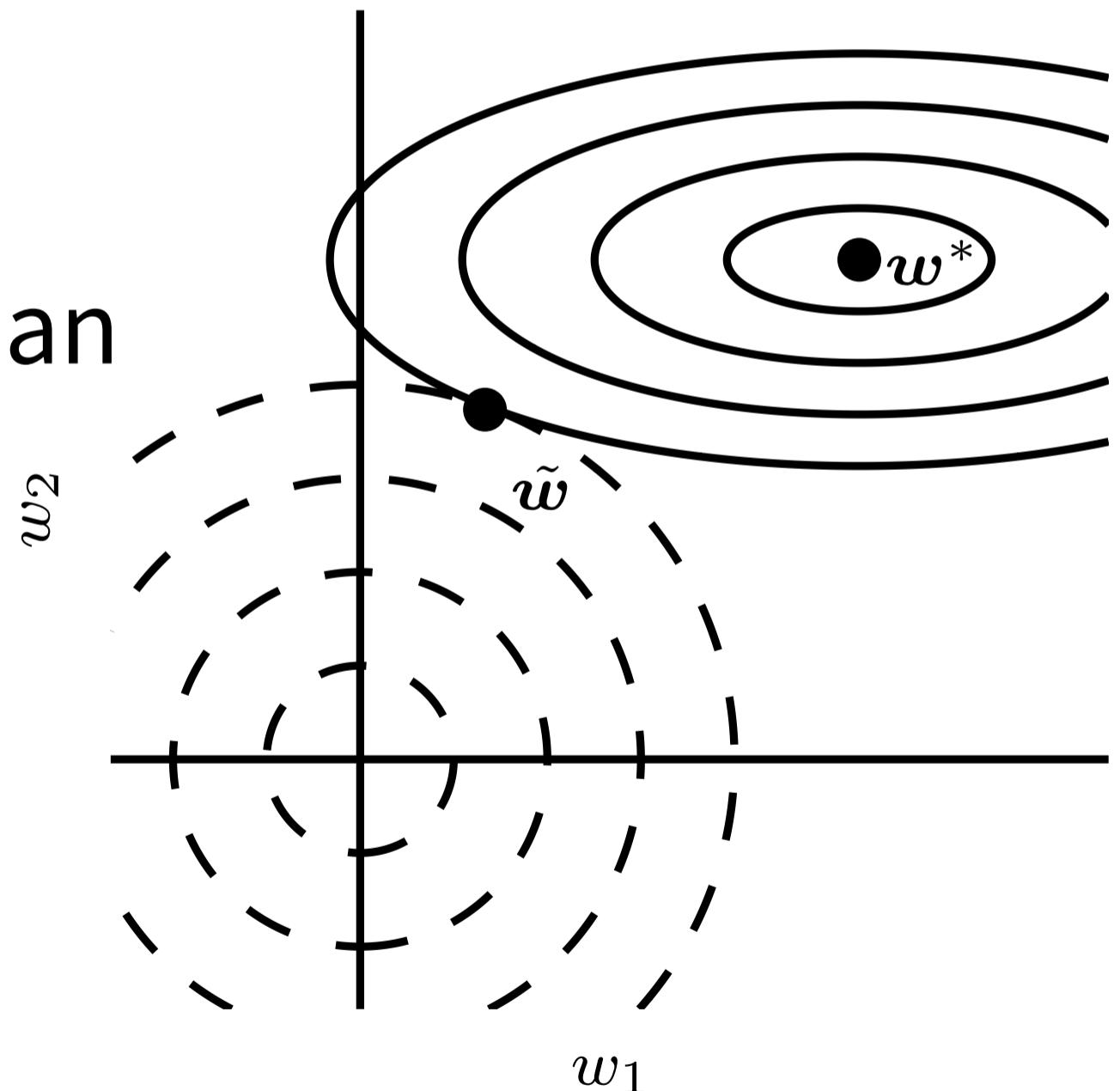
L2 Regularization

Commonly known as **weight decay**:

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

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

Can be interpreted as a Gaussian prior over weights



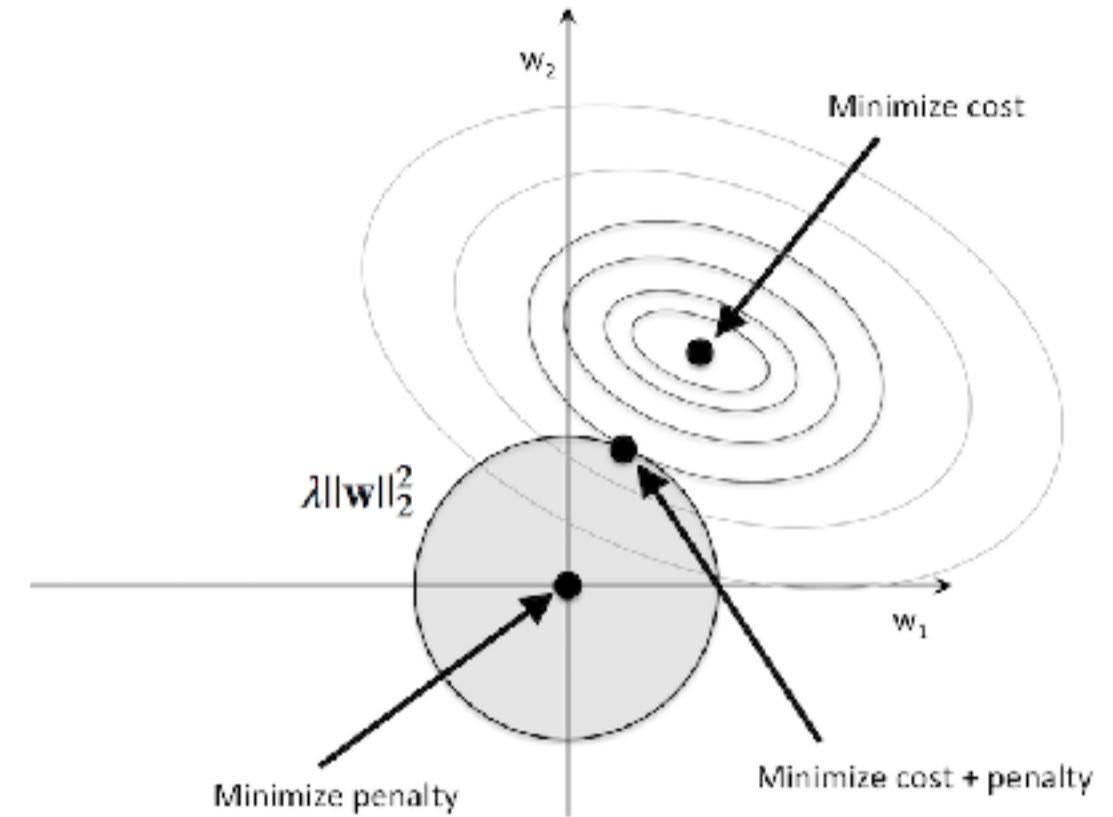
L1 Regularization

L2 regularization

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

Gradient: $\nabla_{W^{(k)}} \Omega(\theta) = \text{sign}(W^{(k)})$

where $\text{sign}(W^{(k)}) = \mathbf{1}_{W_{i,j}^{(k)} > 0} - \mathbf{1}_{W_{i,j}^{(k)} < 0}$

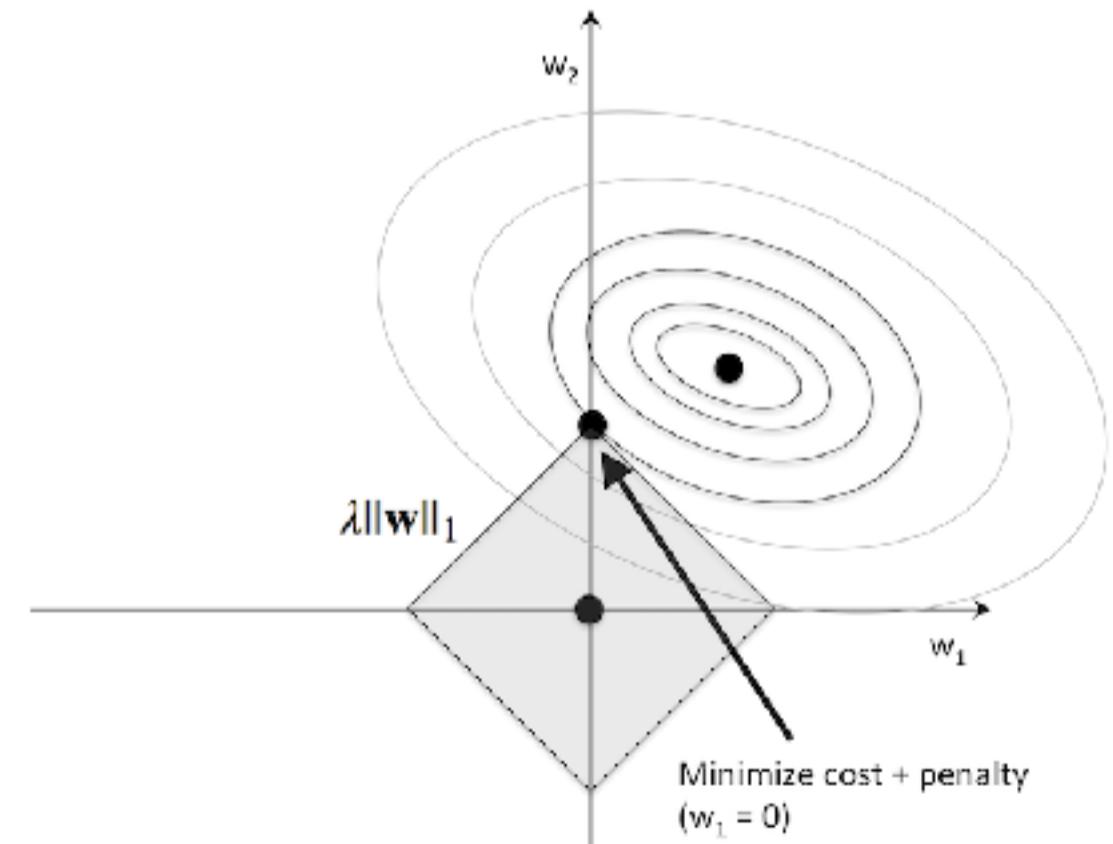


Results in a solution that is more **sparse**

L1 regularization

Will tend to push certain weights to be exactly zero

Can be interpreted as a Gaussian prior over weights



Dataset Augmentation

- Best way to make a ML model generalize better is **collect more data**
- One way to do so is synthesize new data and add it to the training set
- Easiest for classification, in domains where we know how examples should vary (e.g. images)



Affine
Distortion



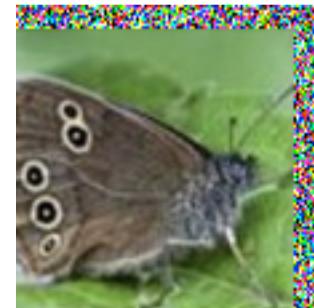
Horizontal
flip



Noise



Random
Translation



Elastic
Deformation



Hue Shift



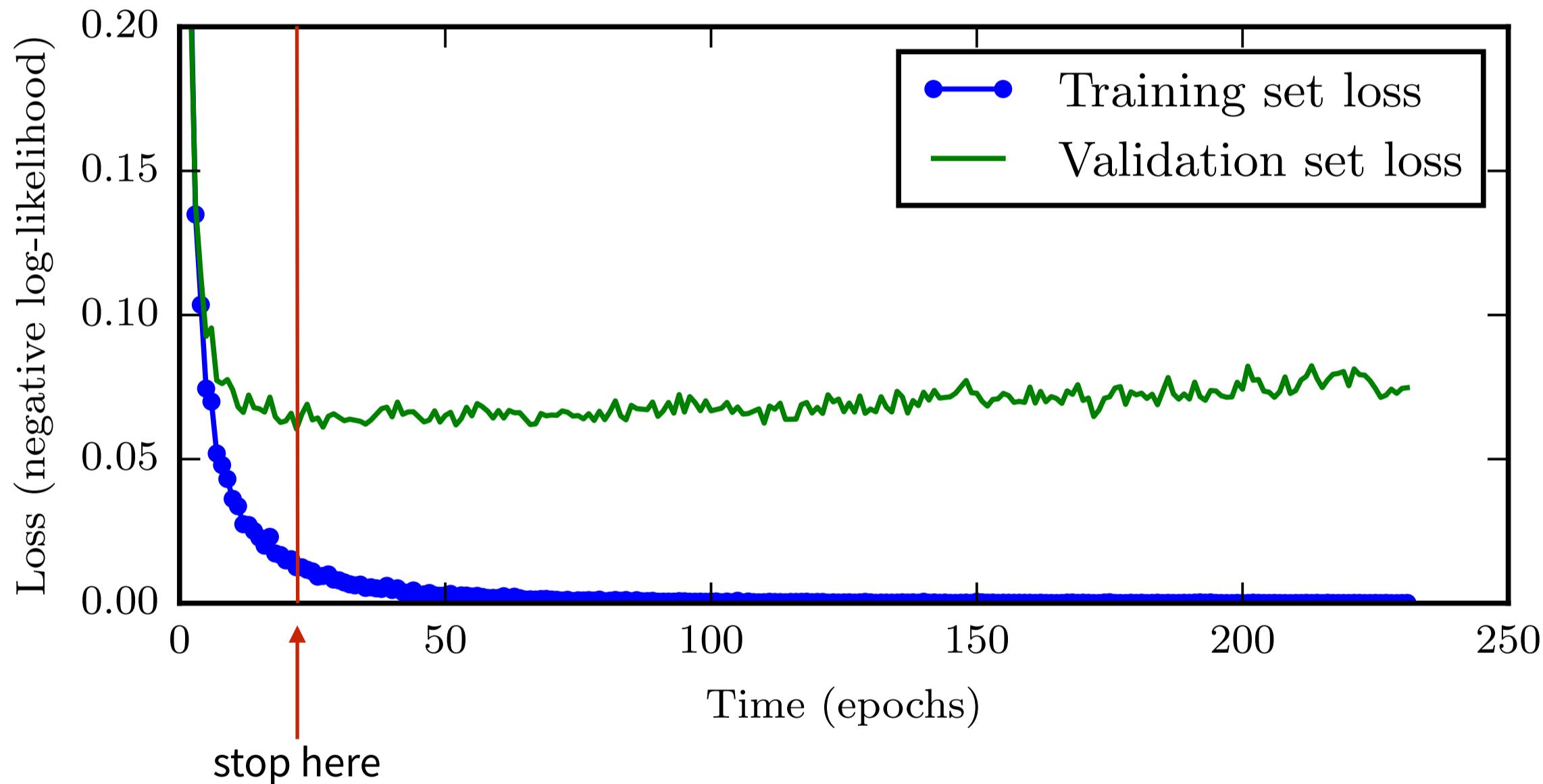
Noise Injection

Various kinds of noise added to networks can encourage robustness:

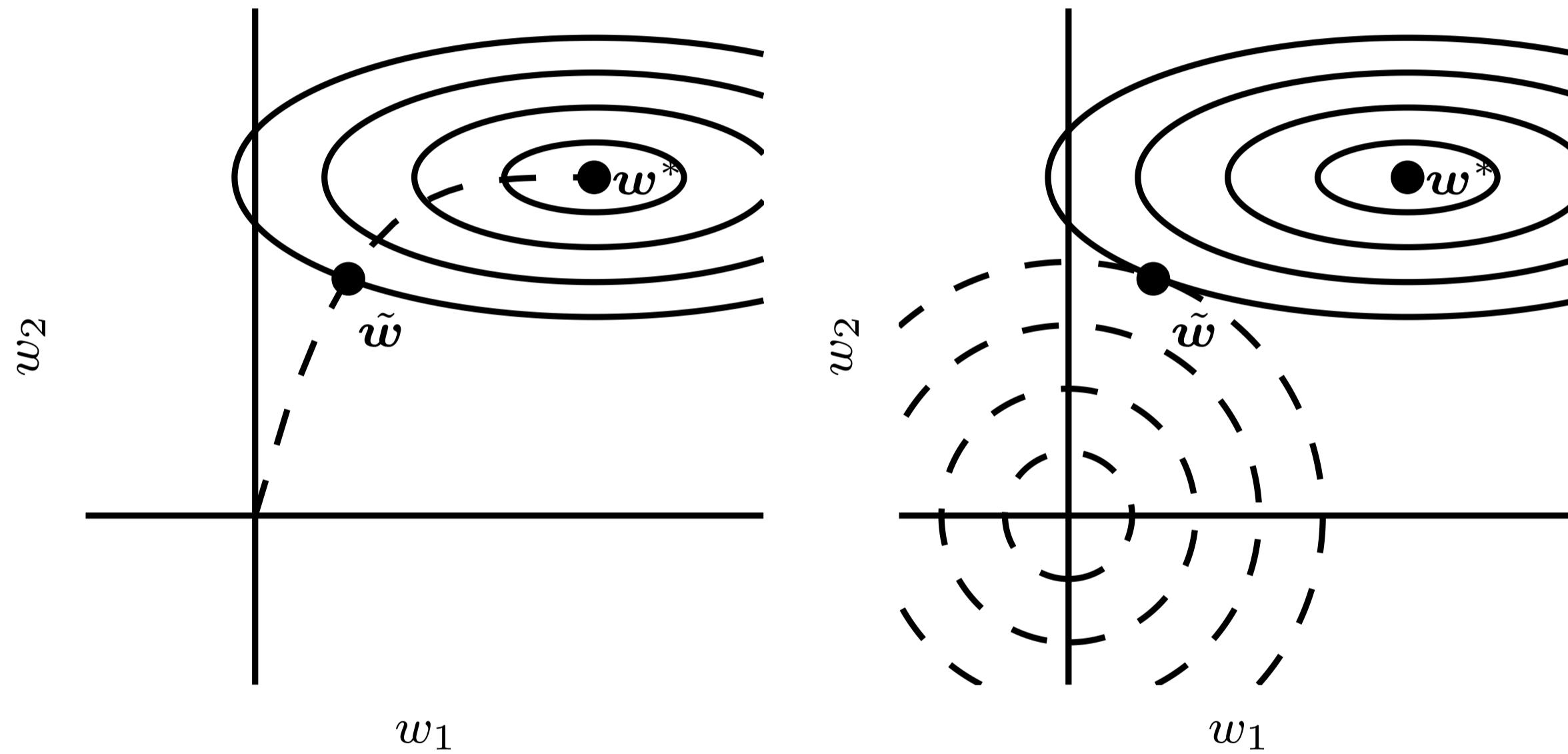
- Noise added to the **inputs**
 - in some models, equivalent to imposing a penalty on the norm of the weights
- Noise added to the **weights**
 - can be interpreted as stochastic Bayesian inference over the weights
- Noise added to the output **targets**
 - **label smoothing** is featured prominently in modern neural networks

Early Stopping

Terminate when validation set performance peaks



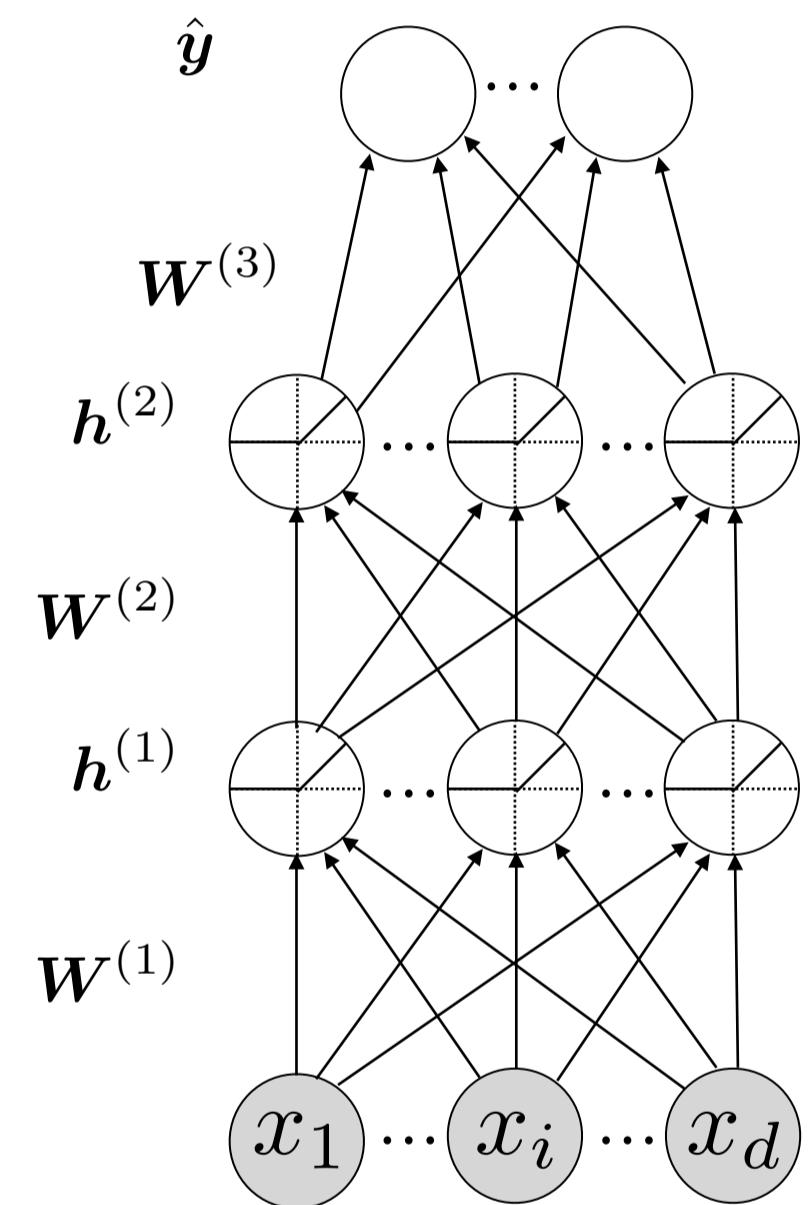
Early Stopping vs. Weight Decay



Dropout

Idea: “cripple” a neural network by removing hidden units stochastically

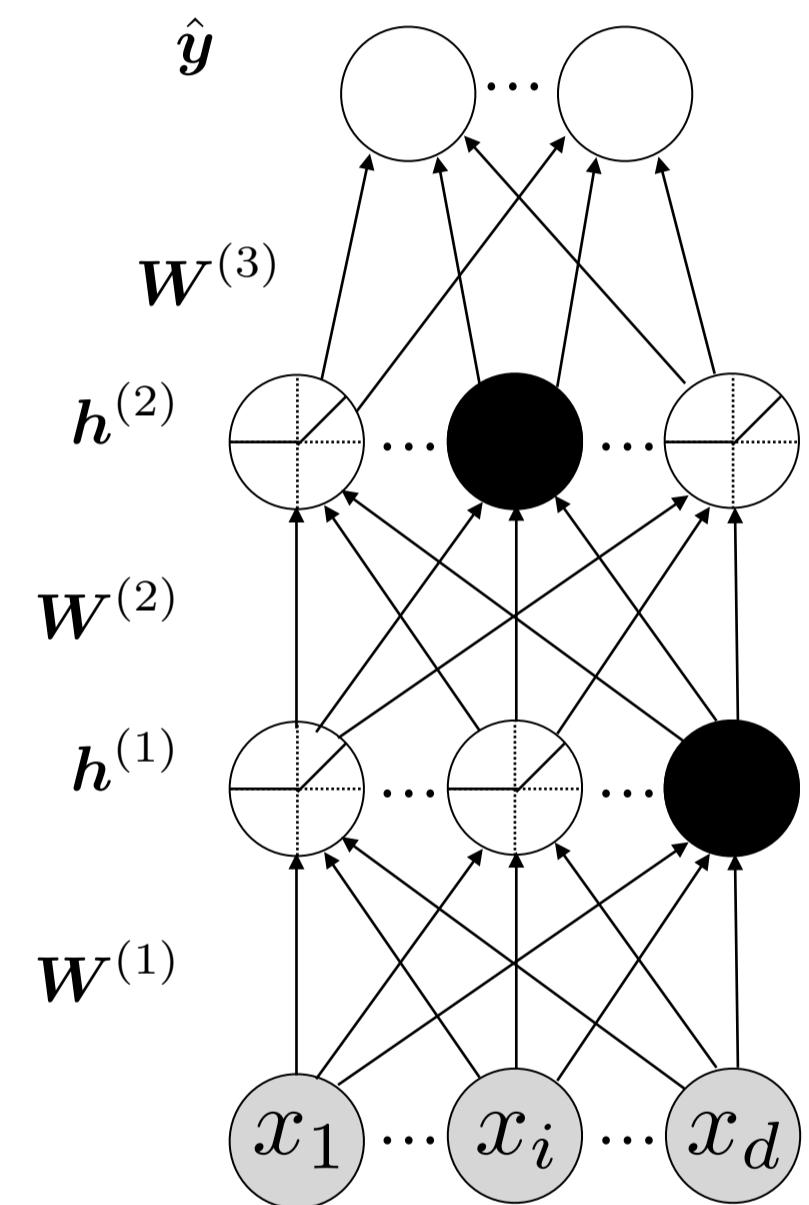
- Set each unit to 0 with probability 0.5 (could use a different probability)
- Hidden units cannot co-adapt to other units
- Causes hidden units to be more useful “on their own”



Dropout

Idea: “cripple” a neural network by removing hidden units stochastically

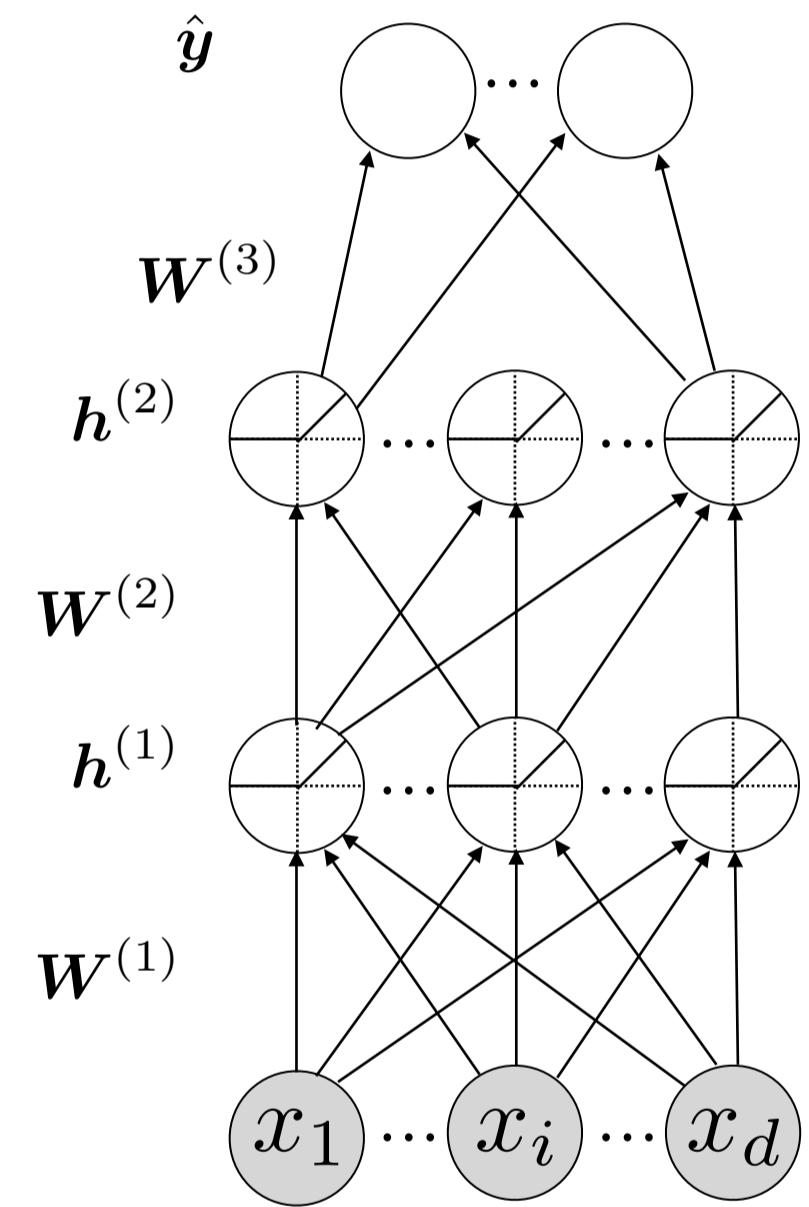
- Set each unit to 0 with probability 0.5 (could use a different probability)
- Hidden units cannot co-adapt to other units
- Causes hidden units to be more useful “on their own”



Dropout: Forward Pass

Implementation: use **random binary masks** $\mathbf{m}^{(k)}$ sampled at training time

- Layer pre-activation for $k > 0$ ($\mathbf{h}^{(0)} = \mathbf{x}$)
$$\mathbf{a}^{(k)} = \mathbf{b}^{(k)} + \mathbf{W}^{(k)} \mathbf{h}^{(k-1)}$$
- Hidden layer activation for $k = 1, \dots, l - 1$
$$\mathbf{h}^{(k)} = g(\mathbf{a}^{(k)}) \odot \mathbf{m}^{(k)}$$
- Output layer activation
$$\hat{\mathbf{y}} = o(\mathbf{a}^{(l)})$$



Dropout: Forward Pass

Implementation: use **random binary masks** $\mathbf{m}^{(k)}$ sampled at training time

- Layer pre-activation for $k > 0$ ($\mathbf{h}^{(0)} = \mathbf{x}$)

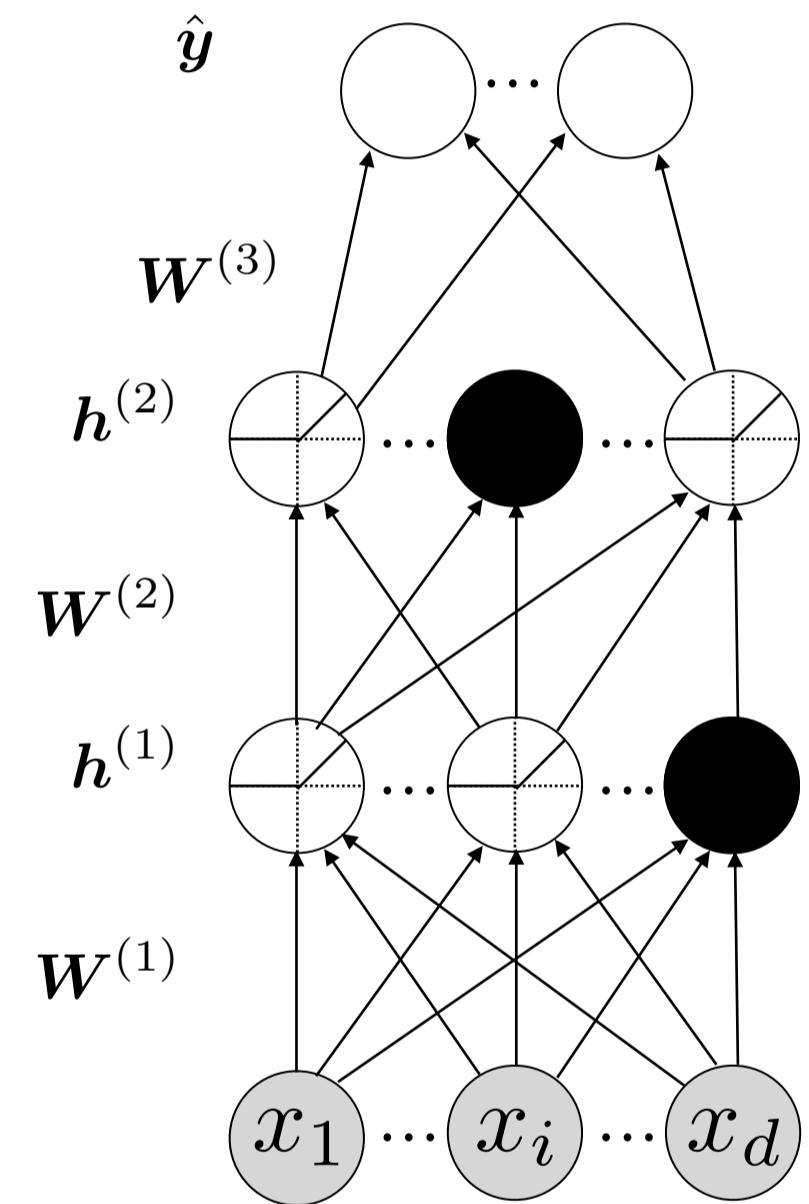
$$\mathbf{a}^{(k)} = \mathbf{b}^{(k)} + \mathbf{W}^{(k)} \mathbf{h}^{(k-1)}$$

- Hidden layer activation for $k = 1, \dots, l - 1$

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

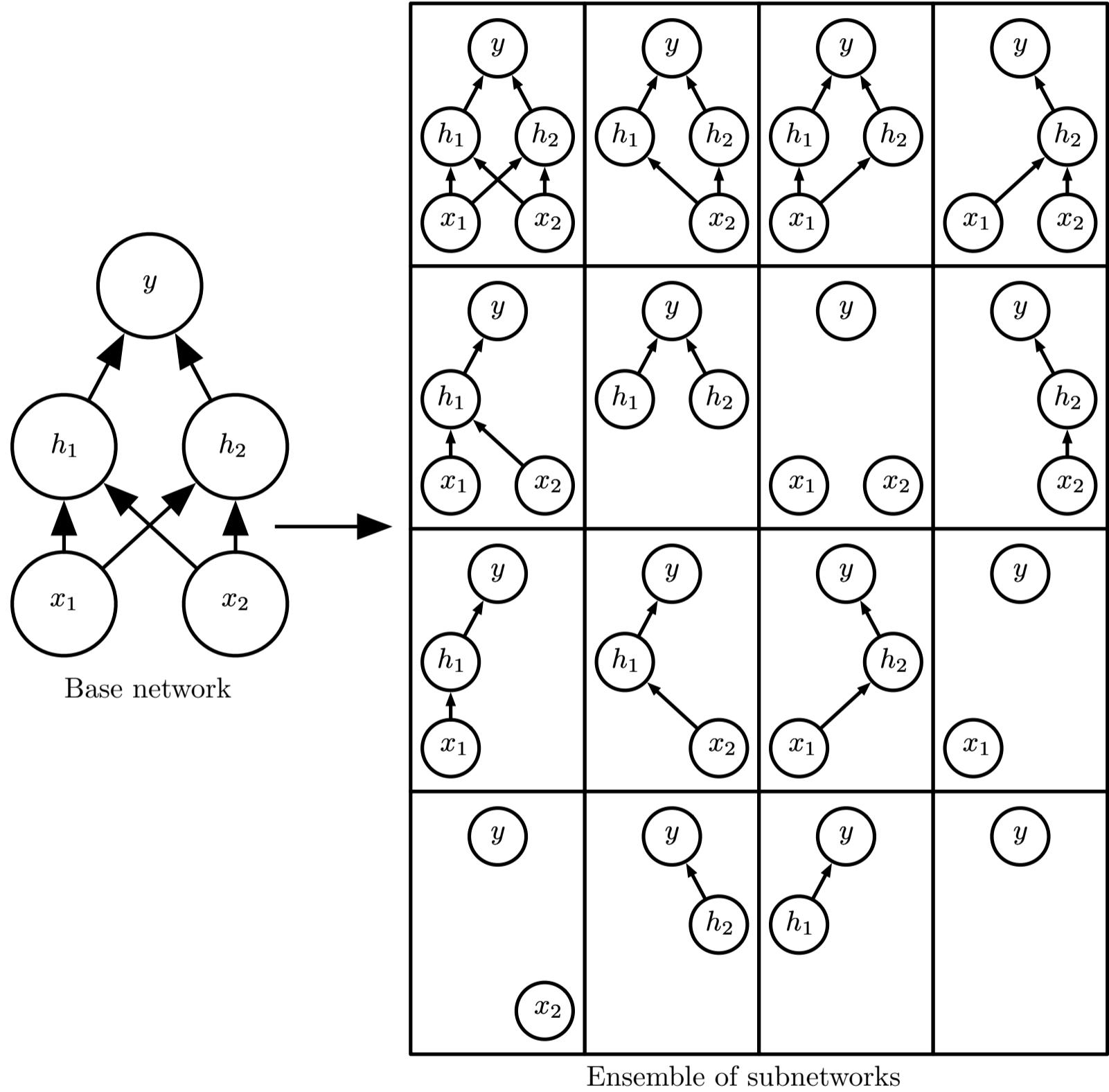
- Output layer activation

$$\hat{\mathbf{y}} = o(\mathbf{a}^{(l)})$$



Dropout: Test Time

- At test time, replace the masks by their expectation
- This is simply the constant vector 0.5 if dropout probability is 0.5
- For a single hidden layer, this is equivalent to taking the **geometric mean** of all neural networks, with all binary masks
- Extremely **efficient ensembling**



Stochastic regularization

Can consider most methods as applying a mask to the weight matrix: $\mathbf{h}^{(k)} = f(\mathbf{W}'^{(k)} \mathbf{h}^{(k-1)}) = f((\mathbf{M}^{(k)} \odot \mathbf{W}^{(k)}) \mathbf{h}^{(k-1)})$

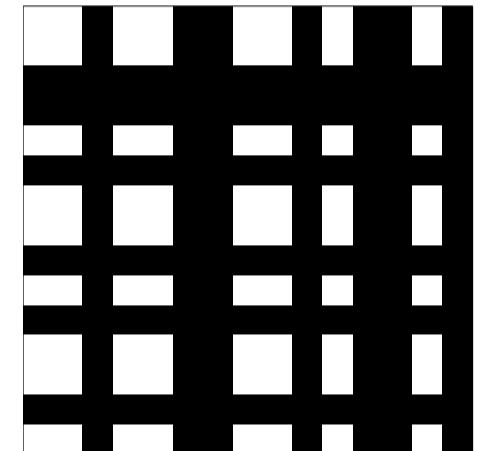
Stochastic regularization

Can consider most methods as applying a mask to the

weight matrix: $\mathbf{h}^{(k)} = f(\mathbf{W}'^{(k)} \mathbf{h}^{(k-1)}) = f((\mathbf{M}^{(k)} \odot \mathbf{W}^{(k)}) \mathbf{h}^{(k-1)})$

Dropout
(Srivastava et al. 2012)

$$\mathbf{M}^{(k)} = \mathbf{m}^{(k)} \mathbf{m}^{(k-1)\top}$$



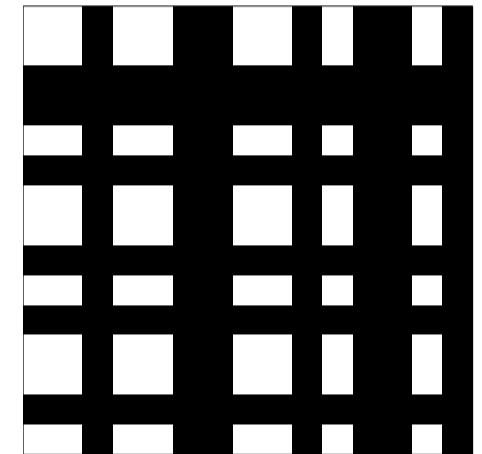
Stochastic regularization

Can consider most methods as applying a mask to the

weight matrix: $\mathbf{h}^{(k)} = f(\mathbf{W}'^{(k)} \mathbf{h}^{(k-1)}) = f((\mathbf{M}^{(k)} \odot \mathbf{W}^{(k)}) \mathbf{h}^{(k-1)})$

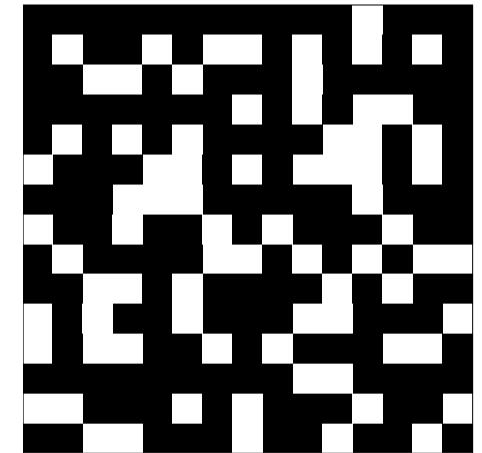
Dropout
(Srivastava et al. 2012)

$$\mathbf{M}^{(k)} = \mathbf{m}^{(k)} \mathbf{m}^{(k-1)\top}$$



DropConnect
(Wan et al. 2013)

$$\mathbf{M}^{(k)} = [\sim \text{Bern}(p^{(k)})]$$



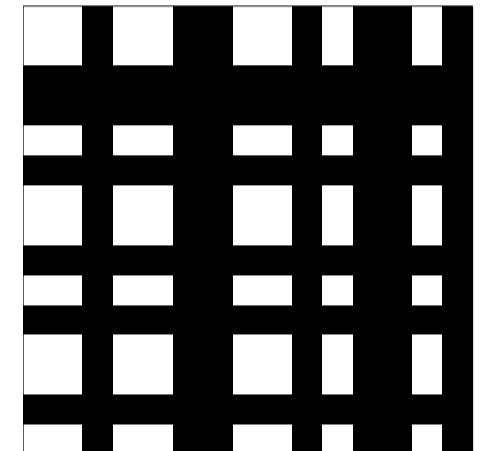
Stochastic regularization

Can consider most methods as applying a mask to the

weight matrix: $\mathbf{h}^{(k)} = f(\mathbf{W}'^{(k)} \mathbf{h}^{(k-1)}) = f((\mathbf{M}^{(k)} \odot \mathbf{W}^{(k)}) \mathbf{h}^{(k-1)})$

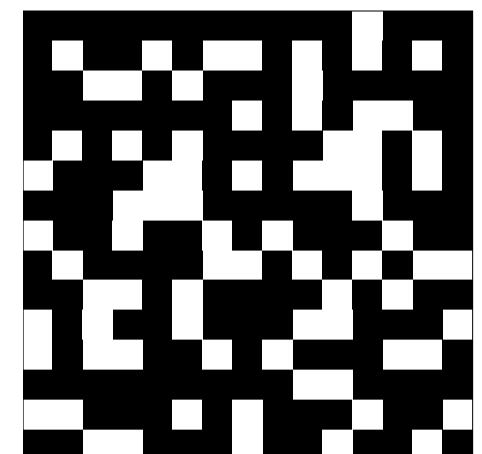
Dropout
(Srivastava et al. 2012)

$$\mathbf{M}^{(k)} = \mathbf{m}^{(k)} \mathbf{m}^{(k-1)\top}$$



DropConnect
(Wan et al. 2013)

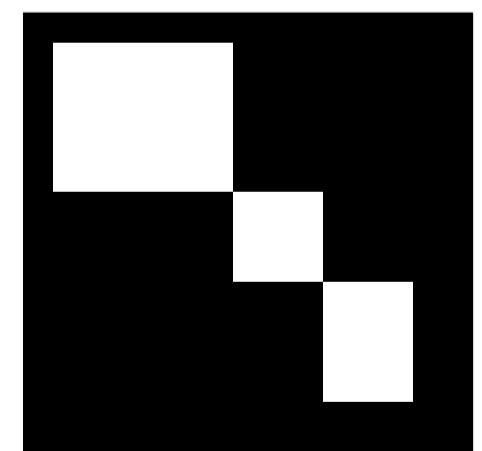
$$\mathbf{M}^{(k)} = [\sim \text{Bern}(p^{(k)})]$$



Blockout
(Murdoch et al. 2016)

$$\begin{aligned}\mathbf{M}^{(k)} &= \frac{1}{c} \mathbf{C}^{(k)} \mathbf{C}^{(k-1)\top} \\ \mathbf{C}^{(k)} &\sim \text{Bern}(\mathbf{P}^{(k)})\end{aligned}$$

c – # of clusters



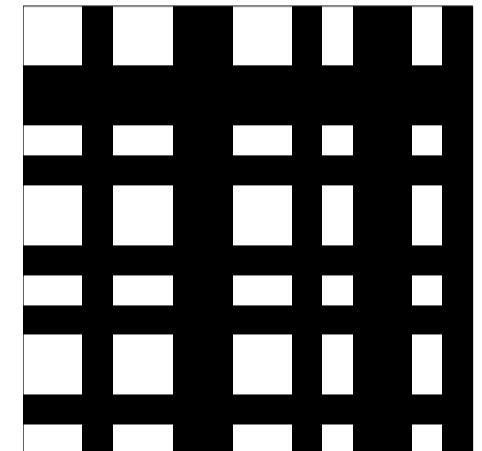
Stochastic regularization

Can consider most methods as applying a mask to the

weight matrix: $\mathbf{h}^{(k)} = f(\mathbf{W}'^{(k)} \mathbf{h}^{(k-1)}) = f((\mathbf{M}^{(k)} \odot \mathbf{W}^{(k)}) \mathbf{h}^{(k-1)})$

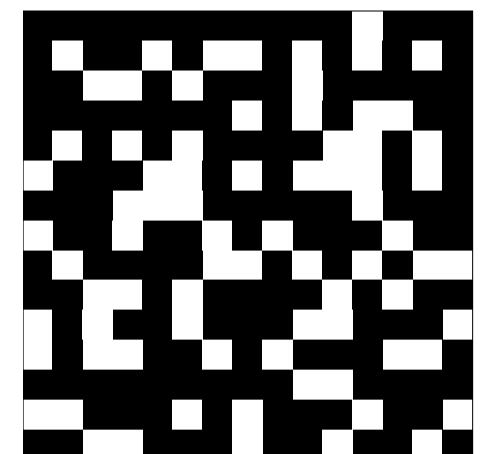
Dropout
(Srivastava et al. 2012)

$$\mathbf{M}^{(k)} = \mathbf{m}^{(k)} \mathbf{m}^{(k-1)\top}$$



DropConnect
(Wan et al. 2013)

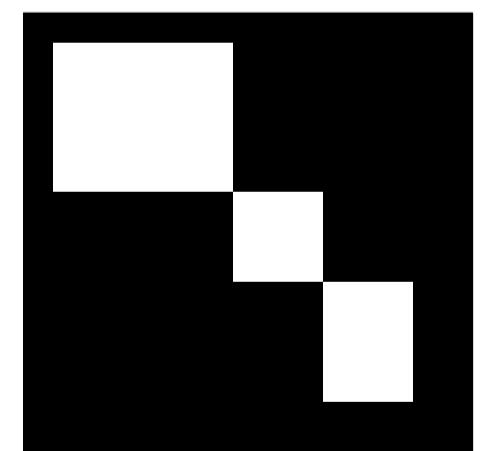
$$\mathbf{M}^{(k)} = [\sim \text{Bern}(p^{(k)})]$$



Blockout
(Murdoch et al. 2016)

$$\begin{aligned}\mathbf{M}^{(k)} &= \frac{1}{c} \mathbf{C}^{(k)} \mathbf{C}^{(k-1)\top} \\ \mathbf{C}^{(k)} &\sim \text{Bern}(\mathbf{P}^{(k)})\end{aligned}$$

c – # of clusters



Blockout: Architecture by Masking

- By stochastically assigning units to clusters, Blockout (Murdoch et al. 2016) permit a “library” of architectural components:

