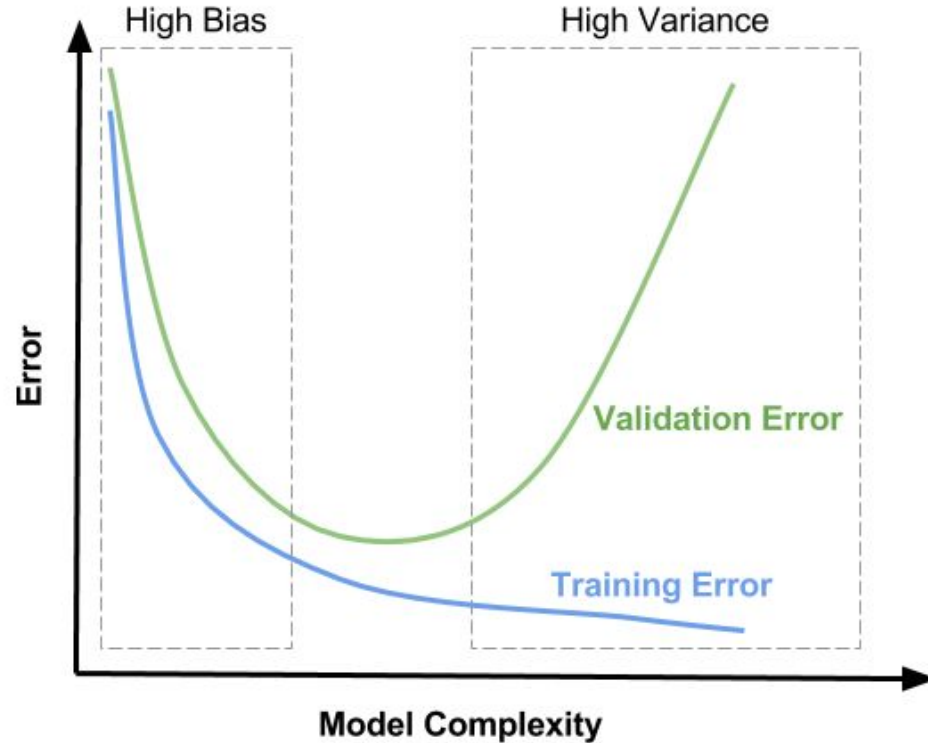


# Neural Networks and Optimization IV

Dr. Parlett-Pelleriti

# Bias Variance Tradeoff



# Regularization

With great complexity comes great regularization

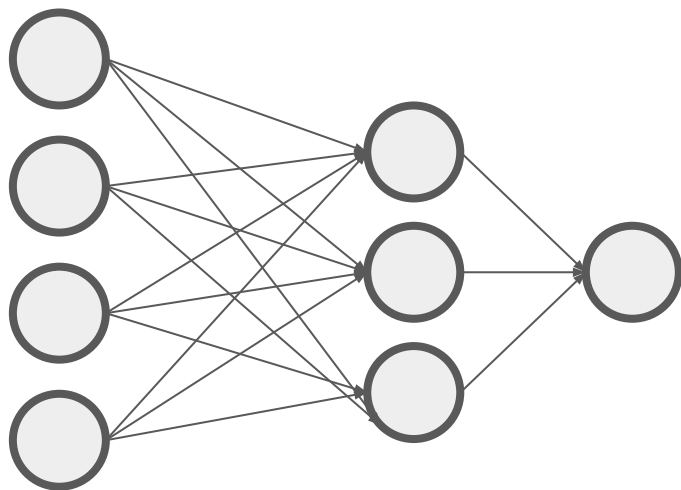
# Regularization

- Penalization
- Bagging
- Dropout
- Early Stopping
- Batch Normalization

# Penalization

- Best for smaller networks
- Added to EACH layer
- L1, L2, L1 + L2

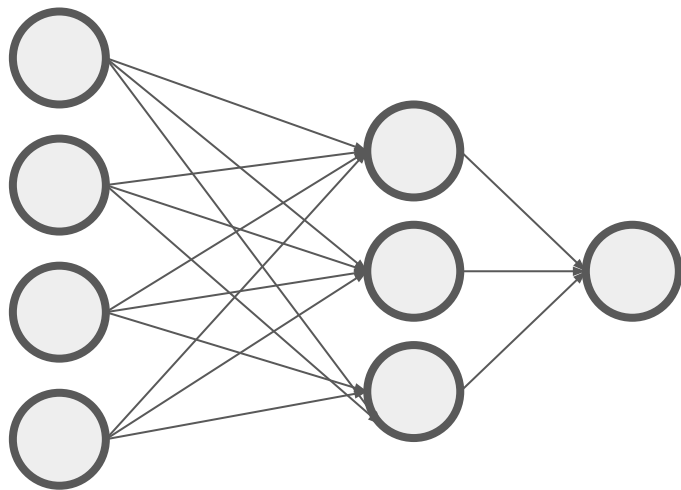
# Penalization



# Penalization

$$\text{penalized loss} = \text{loss} + \lambda \Omega(\theta)$$

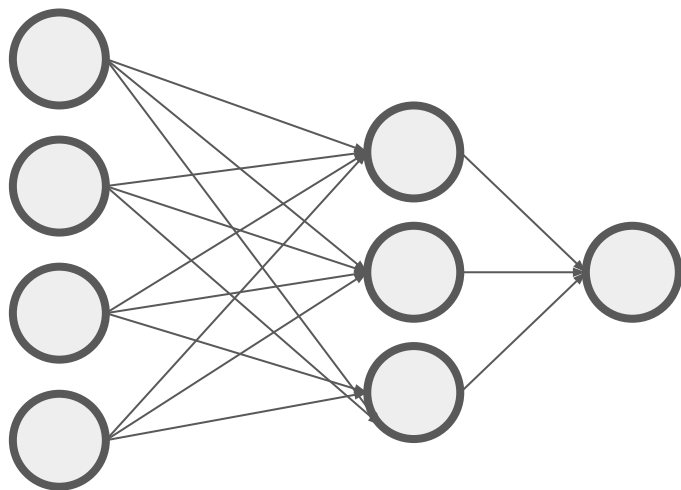
$\Omega(\theta)$  is some penalty on the size of the parameters  $\theta$



# Penalization

$$\text{penalized loss} = \text{loss} + \lambda \Omega(\theta)$$

$\Omega(\theta)$  is some penalty on the size of the parameters  $\theta$



Typically we  
only penalize  
the weights, not  
the biases in a  
network





# Penalization

$$\text{penalized loss} = \text{loss} + \lambda \Omega(\theta)$$

$\Omega(\theta)$  is some penalty on the size of the parameters  $\theta$

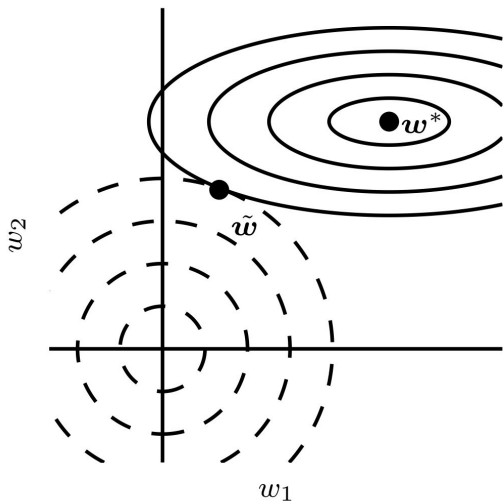


Figure 7.1

## L<sub>2</sub> Penalization

$$\text{penalized loss} = \text{loss} + \lambda \Omega(\theta)$$

$\Omega(\theta)$  is some penalty on the size of the parameters  $\theta$

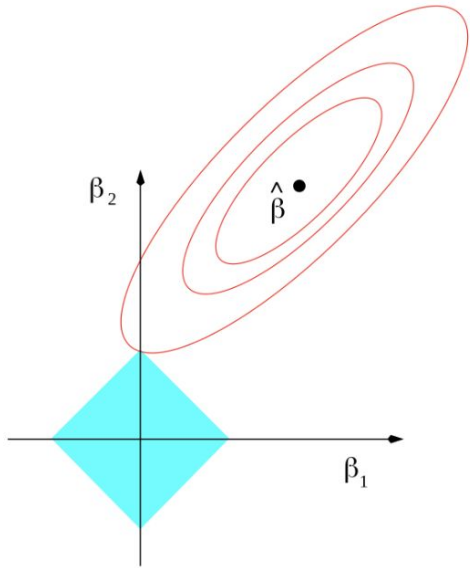
$$w_t = \underbrace{w_{t-1}}_{\text{old weights}} - \underbrace{\alpha * g_t}_{\text{update weights}}$$

$$w_t = \underbrace{(1 - \lambda\alpha)w_{t-1}}_{\text{shrink weights}} - \underbrace{\alpha * g_t}_{\text{update weights}}$$

# $L_1$ Penalization

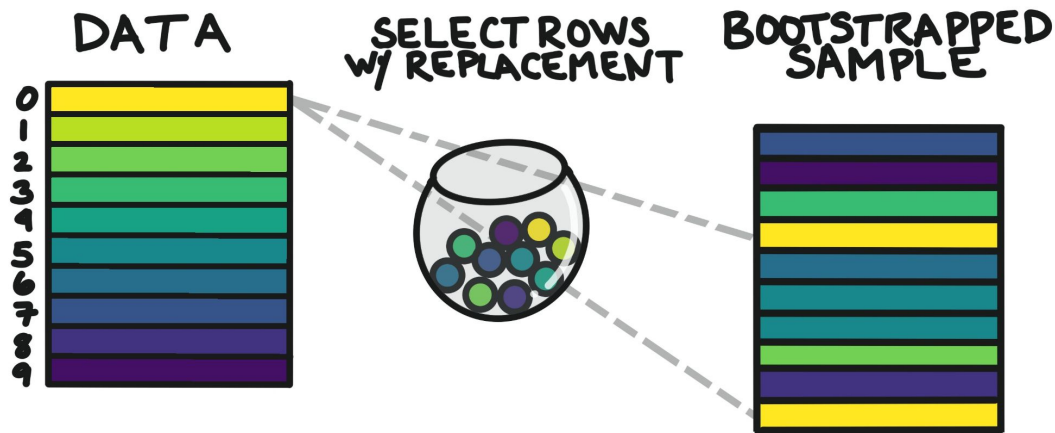
$$\text{penalized loss} = \text{loss} + \lambda \Omega(\theta)$$

$\Omega(\theta)$  is some penalty on the size of the parameters  $\theta$



# Bagging

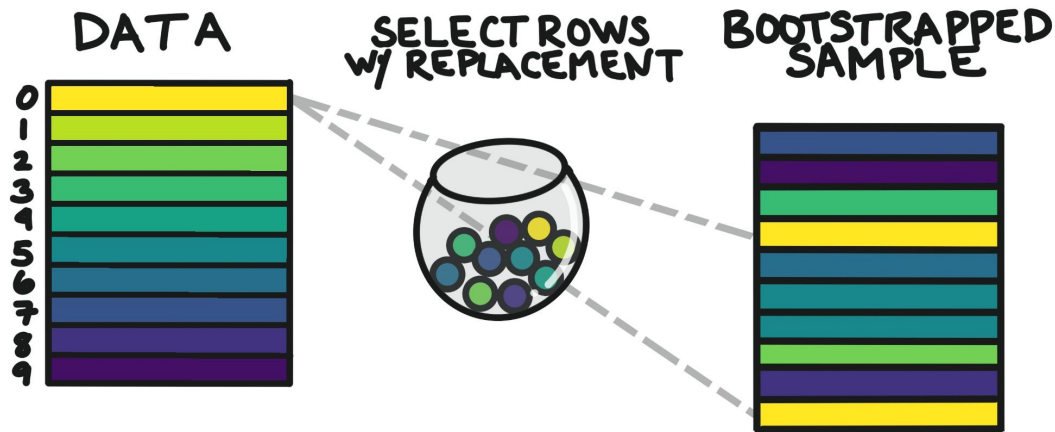
## BOOTSTRAPPING



@CHELSEAPARLETT

# Bagging

## BOOTSTRAPPING

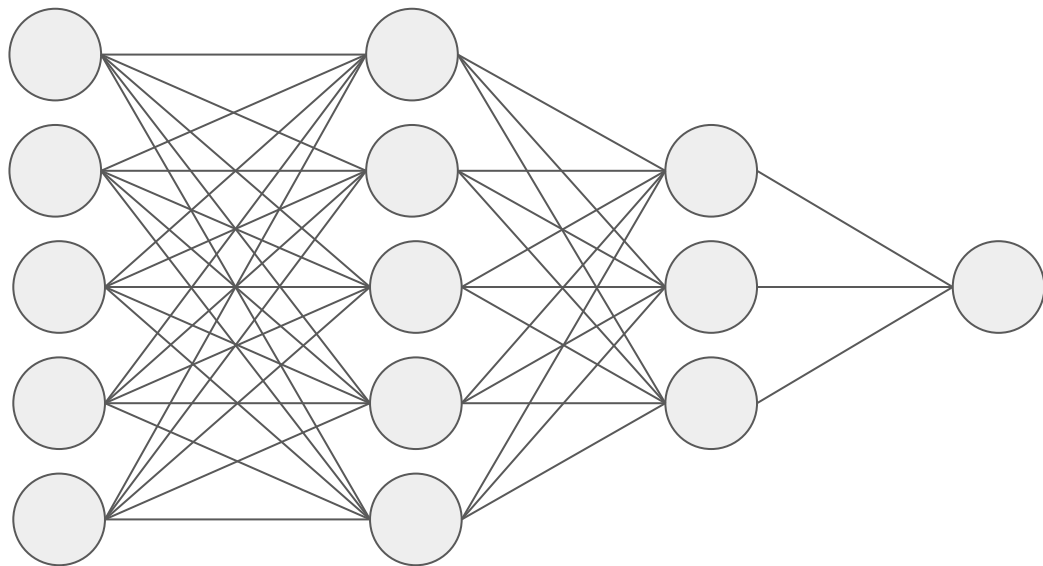


@CHELSEAPARLETT

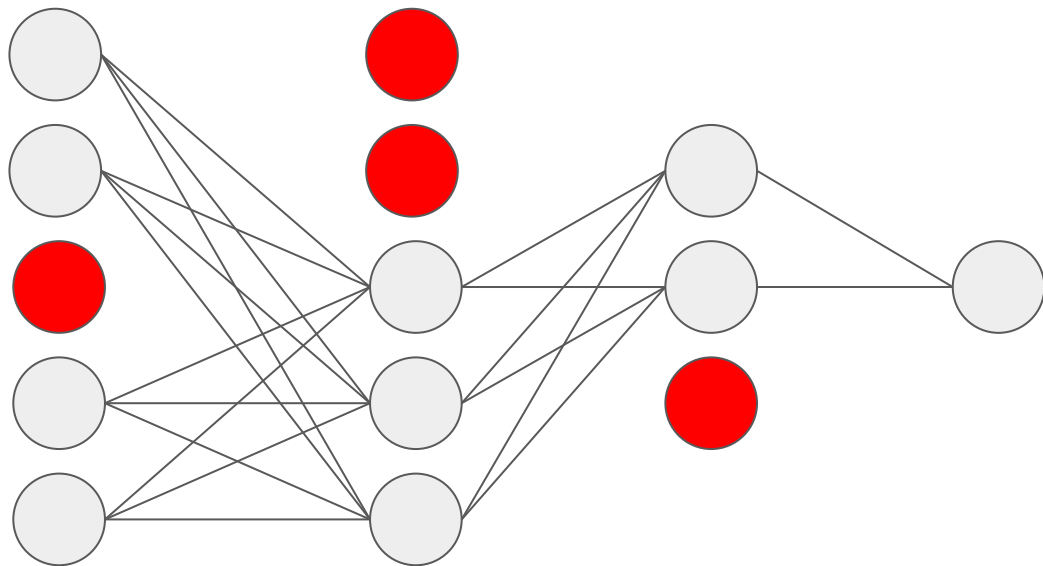
# Dropout

- Comparison to Random Forest Bagging/Feature Selection
- Applied to each layer
- Network can't over-rely on some connections

# Dropout

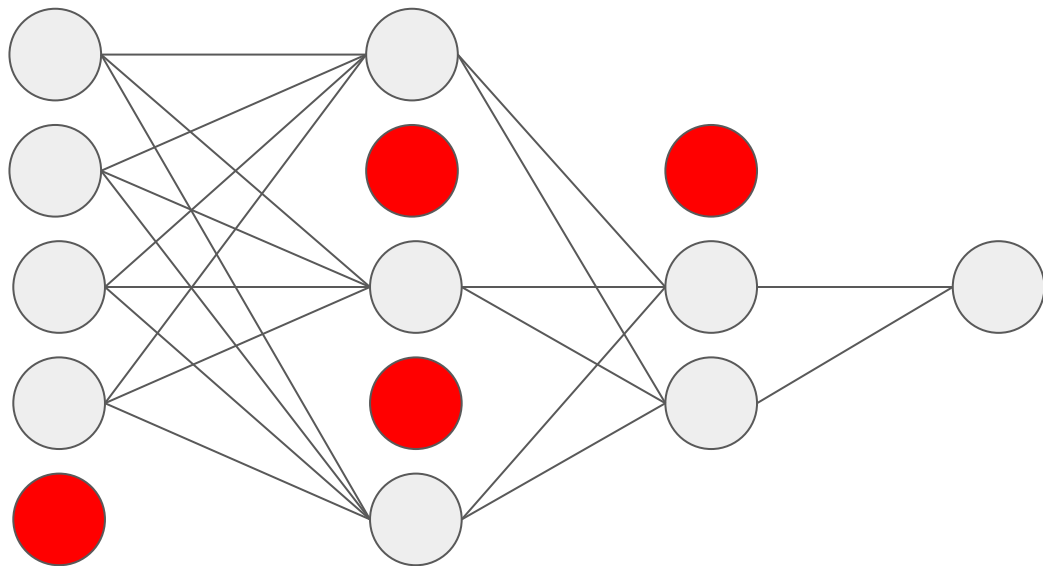


# Dropout

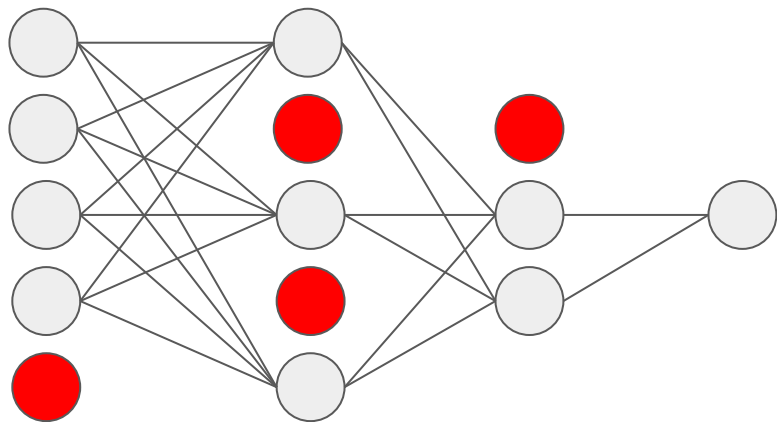




# Dropout

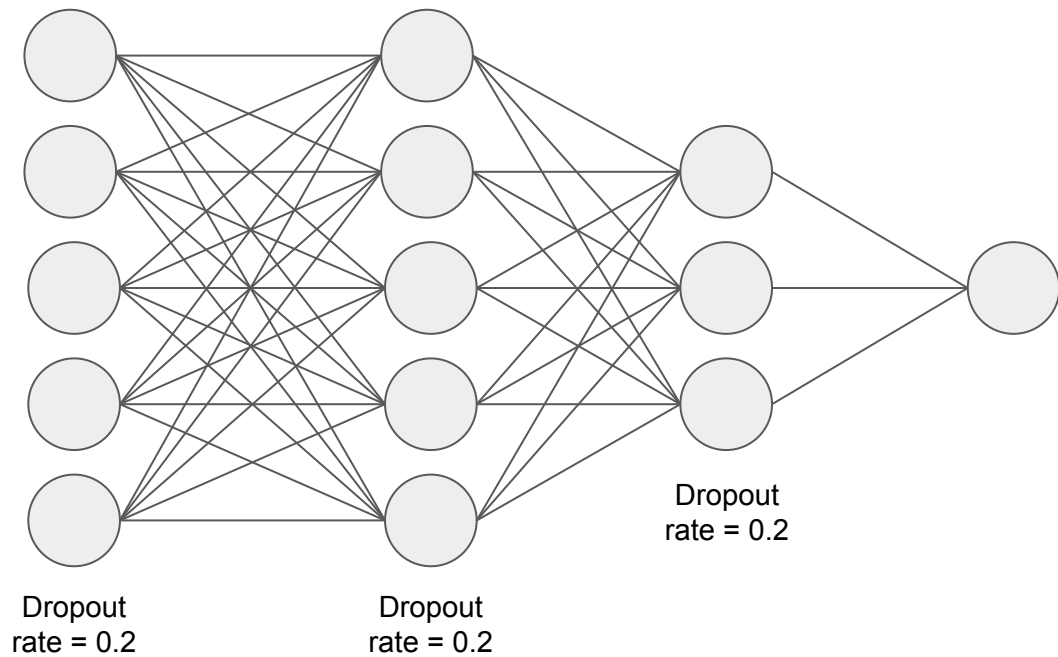


# Dropout

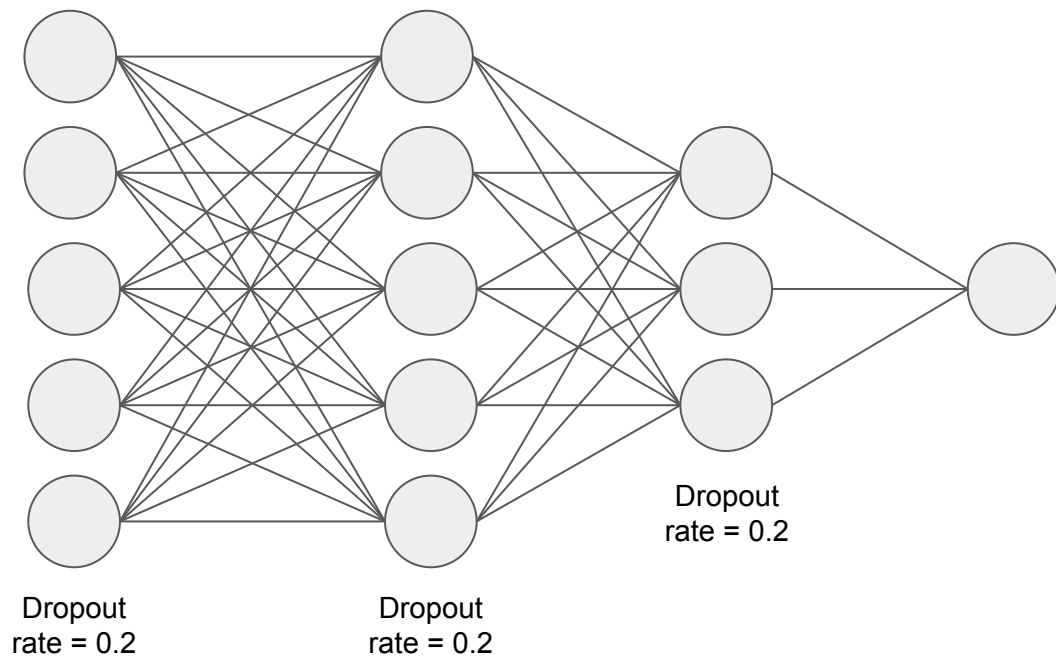


|       |     |      |       |      |
|-------|-----|------|-------|------|
| 0.1   | 0.4 | 0.2  | 0.3   | 0.2  |
| 0.1   | 0   | 0.2  | 0.3   | 0.2  |
| 0.125 | 0   | 0.25 | 0.375 | 0.25 |

# Dropout

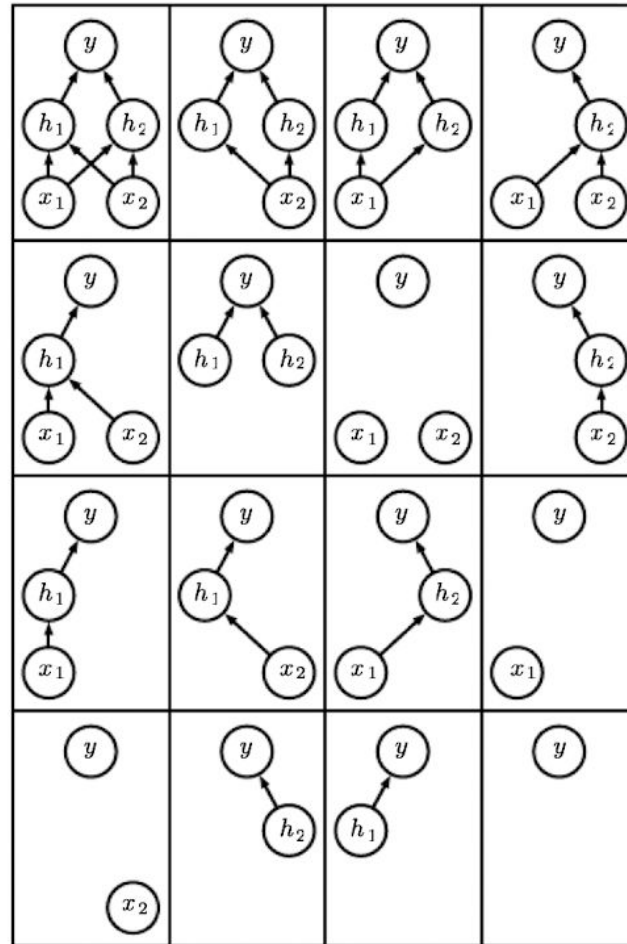
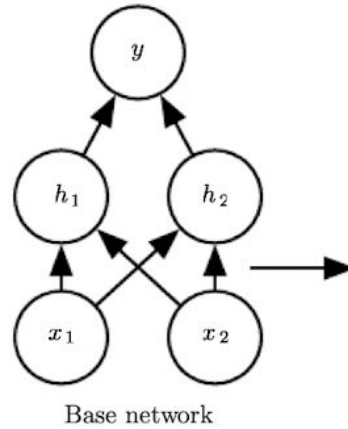


# Dropout



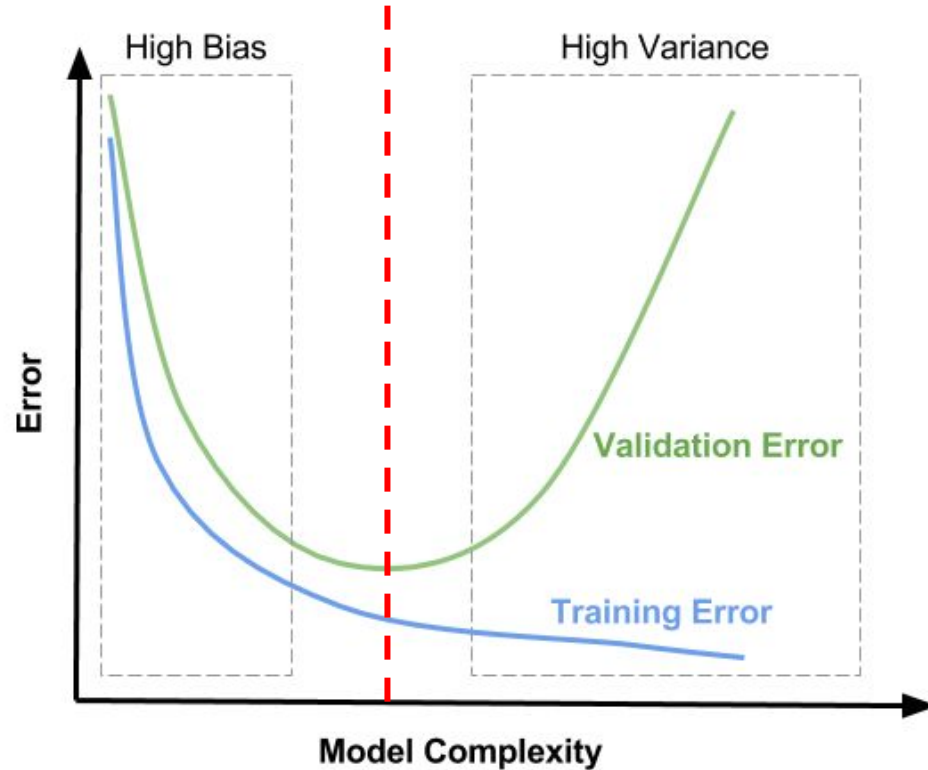
**Dropout Only  
Happens During  
Training!**

# Dropout

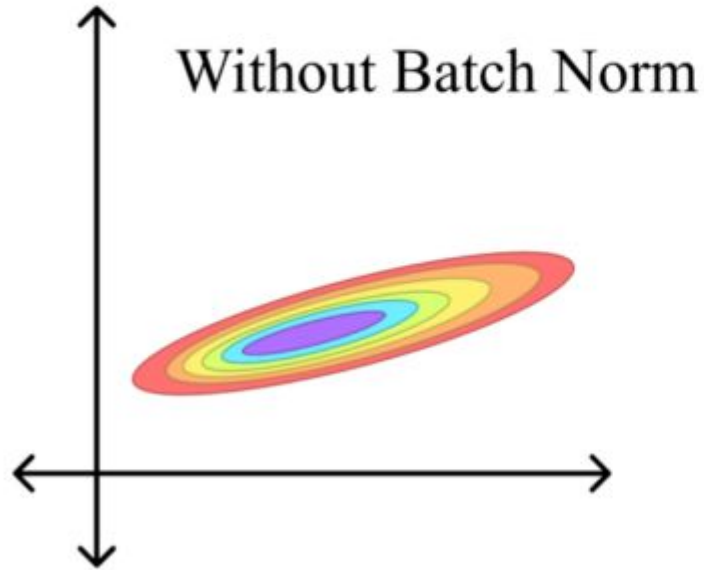


Ensemble of subnetworks

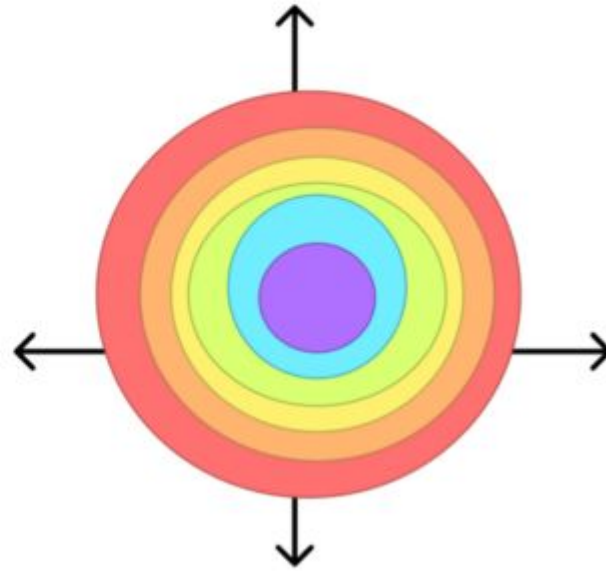
# Early Stopping



# Batch Normalization



With Batch Norm



# Batch Normalization

- Speeds up training
- Lessens the impact of initial weights
- Regularizes your model



**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_{1...m}\}$ ;

Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma,\beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma,\beta}(x_i) \quad // \text{ scale and shift}$$

**Algorithm 1:** Batch Normalizing Transform, applied to activation  $x$  over a mini-batch.