

# Regularization

Goodfellow et al. 2016 (sec. 7)

# Regularization: linear regression

- Recall the linear regression problem

$$\hat{L}(\mathbf{w}) = \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \quad \mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- Inverse well-defined only when  $\text{rank}(\mathbf{X}^T \mathbf{X}) = d$
- Add **regularization term**

$$\hat{L}(\mathbf{w}) = \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \beta \|\mathbf{w}\|^2$$

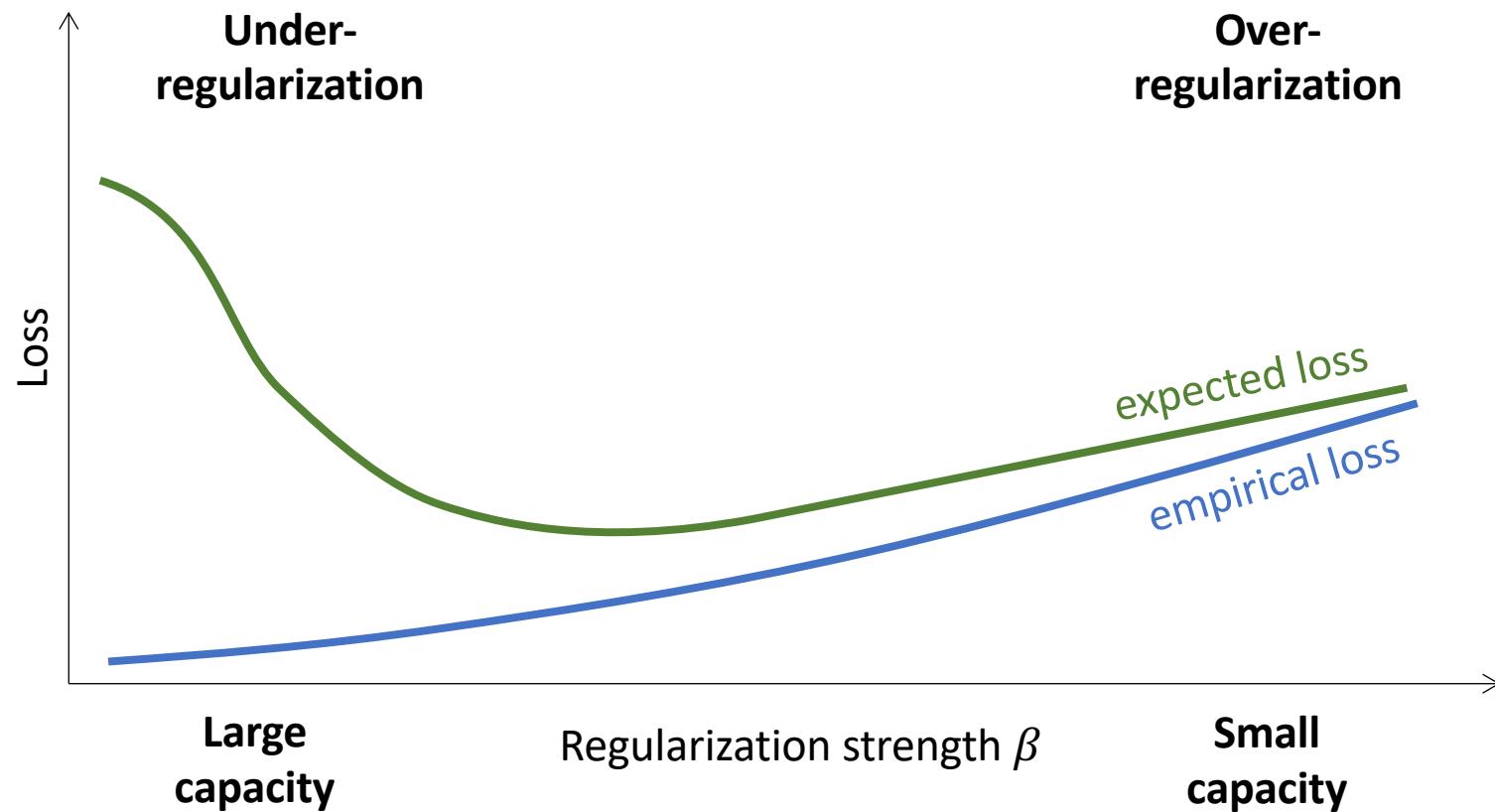
$$\mathbf{w}^* = \left( \frac{1}{n} \mathbf{X}^T \mathbf{X} + \beta \mathbf{I} \right)^{-1} \mathbf{X}^T \mathbf{y}$$

# Regularization: general form

$$\min_{\theta} \hat{L}(\theta) + \beta R(\theta)$$

- Soft constraint (“penalty”)
- Equivalent to hard constraint (Lagrangian)
- Maximum a posteriori interpretation: principled way of deriving  $R$  if we know the distribution of  $\theta$
- Restricts the hypothesis class
- Helps prevent overfitting

# Regularization



# Bayesian view

$$p(\mathbf{x}, y | \boldsymbol{\theta}) = \frac{p(\mathbf{x}, y) p(\boldsymbol{\theta} | \mathbf{x}, y)}{p(\boldsymbol{\theta})}$$

- Prior on hypotheses (parameters distribution):  $p(\boldsymbol{\theta})$

- Maximum a posteriori (MAP):

$$\min_{\boldsymbol{\theta}} - \log p(\boldsymbol{\theta} | \mathbf{x}, y) = \min_{\boldsymbol{\theta}} - \underbrace{\log p(\boldsymbol{\theta})}_{\text{Regularization}} - \underbrace{\log p(\mathbf{x}, y | \boldsymbol{\theta})}_{\text{MLE}}$$

# $L_2$ Regularization

$$\hat{L}_R(\boldsymbol{\theta}) = \hat{L}(\boldsymbol{\theta}) + \frac{\beta}{2} \|\boldsymbol{\theta}\|^2$$

- Gradient of regularized loss becomes

$$\nabla_{\boldsymbol{\theta}} \hat{L}_R(\boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \hat{L}(\boldsymbol{\theta}) + \beta \boldsymbol{\theta}$$

- Effect on gradient descent update

$$\begin{aligned}\boldsymbol{\theta}^{(t+1)} &= \boldsymbol{\theta}^{(t)} - \alpha^{(t)} \nabla_{\boldsymbol{\theta}} \hat{L}(\boldsymbol{\theta}^{(t)}) - \alpha^{(t)} \beta \boldsymbol{\theta}^{(t)} \\ &= (1 - \alpha^{(t)} \beta) \boldsymbol{\theta}^{(t)} - \alpha^{(t)} \nabla_{\boldsymbol{\theta}} \hat{L}(\boldsymbol{\theta}^{(t)})\end{aligned}$$

“weight decay”

- Bayesian interpretation:  $\boldsymbol{\theta} \sim \mathcal{N}\left(0, \frac{1}{\sqrt{2\beta}} \mathbf{I}\right)$

# Geometric interpretation

- Consider second-order Taylor expansion around original problem optimum  $\theta^*$

$$\begin{aligned}\hat{L}(\theta) \approx \hat{L}(\theta^*) + \cancel{\nabla_{\theta} \hat{L}(\theta^*)^T (\theta - \theta^*)} \\ + \frac{1}{2} (\theta - \theta^*)^T \nabla_{\theta}^2 \hat{L}(\theta^*) (\theta - \theta^*)\end{aligned}$$

# Geometric interpretation

- Consider second-order Taylor expansion around original problem optimum  $\boldsymbol{\theta}^*$

$$\hat{L}(\boldsymbol{\theta}) \approx \hat{L}(\boldsymbol{\theta}^*) + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)^T \nabla_{\boldsymbol{\theta}}^2 \hat{L}(\boldsymbol{\theta}^*) (\boldsymbol{\theta} - \boldsymbol{\theta}^*)$$

- Gradient of the loss:

$$\nabla_{\boldsymbol{\theta}} \hat{L}(\boldsymbol{\theta}) \approx \nabla_{\boldsymbol{\theta}}^2 \hat{L}(\boldsymbol{\theta}^*) (\boldsymbol{\theta} - \boldsymbol{\theta}^*)$$

# Geometric interpretation

- Consider second-order Taylor expansion around original problem optimum  $\boldsymbol{\theta}^*$

$$\hat{L}(\boldsymbol{\theta}) \approx \hat{L}(\boldsymbol{\theta}^*) + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}^*)^T \nabla_{\boldsymbol{\theta}}^2 \hat{L}(\boldsymbol{\theta}^*) (\boldsymbol{\theta} - \boldsymbol{\theta}^*)$$

- Gradient of the regularized loss:

$$\nabla_{\boldsymbol{\theta}} \hat{L}_R(\boldsymbol{\theta}) \approx \nabla_{\boldsymbol{\theta}}^2 \hat{L}(\boldsymbol{\theta}^*) (\boldsymbol{\theta} - \boldsymbol{\theta}^*) + \beta \boldsymbol{\theta}$$

- Optimality condition:

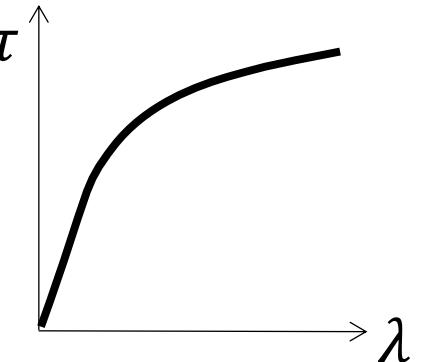
$$\mathbf{0} = \nabla_{\boldsymbol{\theta}} \hat{L}_R(\boldsymbol{\theta}_R^*) \approx \nabla_{\boldsymbol{\theta}}^2 \hat{L}(\boldsymbol{\theta}^*) (\boldsymbol{\theta} - \boldsymbol{\theta}^*) + \beta \boldsymbol{\theta}_R^*$$

$$\boldsymbol{\theta}_R^* \approx (\nabla_{\boldsymbol{\theta}}^2 \hat{L}(\boldsymbol{\theta}^*) + \beta \mathbf{I})^{-1} \nabla_{\boldsymbol{\theta}}^2 \hat{L}(\boldsymbol{\theta}^*) \boldsymbol{\theta}^*$$

# Geometric interpretation

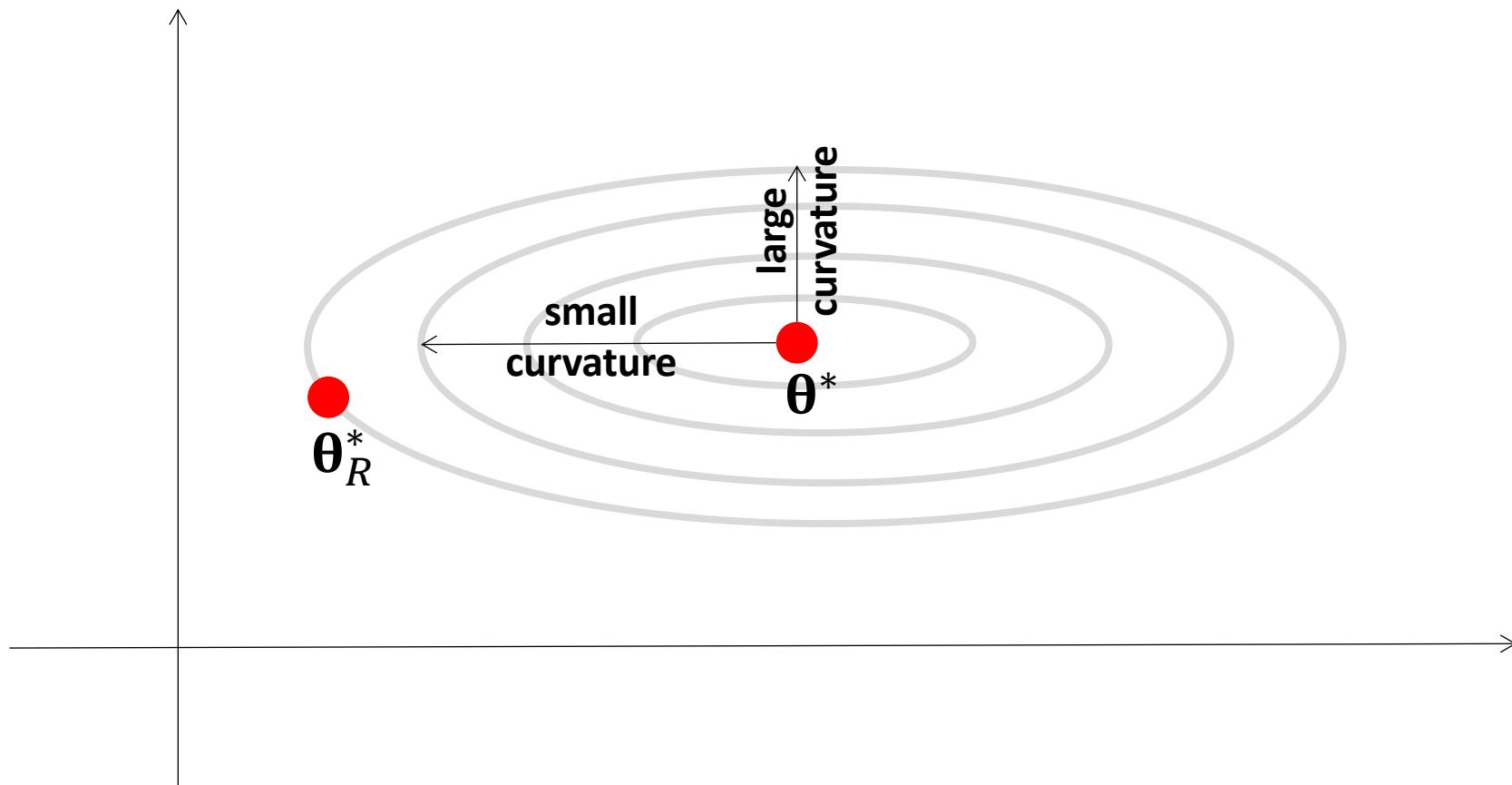
- Assume Hessian has eigendecomposition  $\nabla_{\theta}^2 \hat{L}(\theta^*) = \mathbf{U} \Lambda \mathbf{U}^T$
- Then

$$\begin{aligned}\theta_R^* &\approx (\nabla_{\theta}^2 \hat{L}(\theta^*) + \beta \mathbf{I})^{-1} \nabla_{\theta}^2 \hat{L}(\theta^*) \theta^* \\ &\approx \mathbf{U} (\Lambda + \beta \mathbf{I})^{-1} \Lambda \mathbf{U}^T \theta^* \\ &= \mathbf{U} \tau(\Lambda) \mathbf{U}^T \theta^* \\ &= \tau \left( \nabla_{\theta}^2 \hat{L}(\theta^*) \right) \theta^*\end{aligned}$$



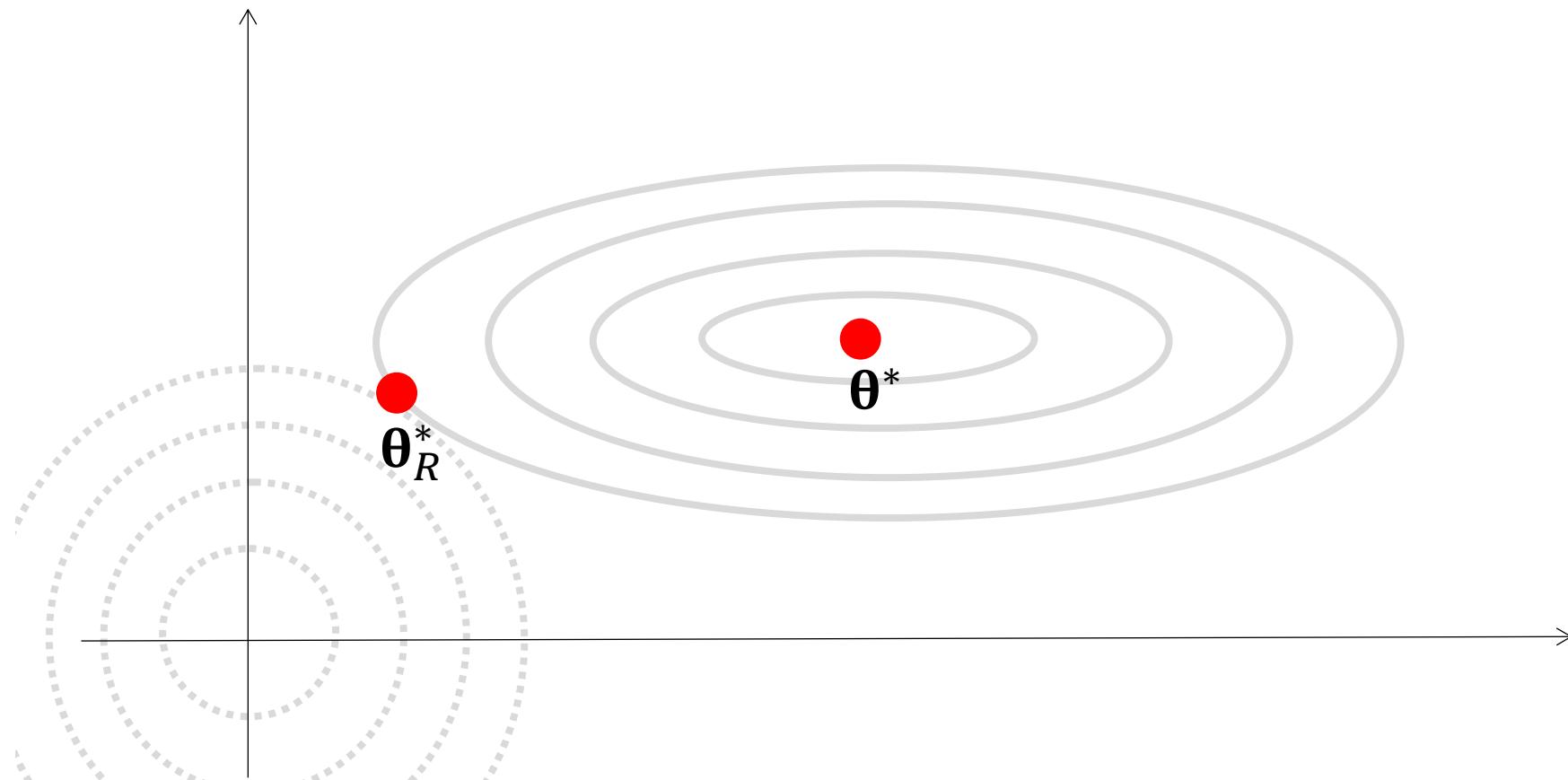
- **Filter of  $\theta^*$**  (attenuates along small eigenvectors)

# Geometric interpretation



Adapted from Goodfellow et al. 2016

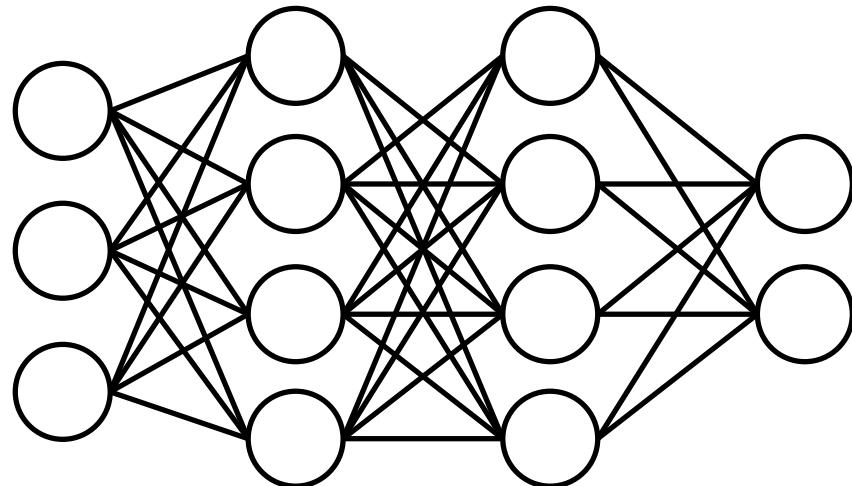
# Geometric interpretation



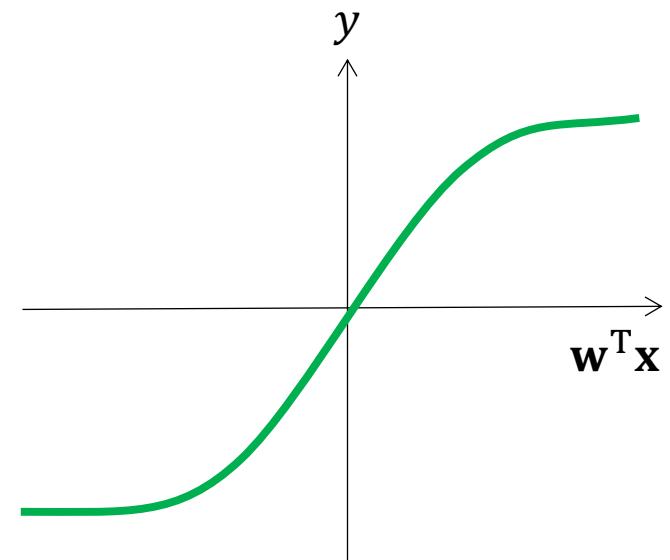
Adapted from Goodfellow et al. 2016

# Intuition

Why small weights reduce capacity?



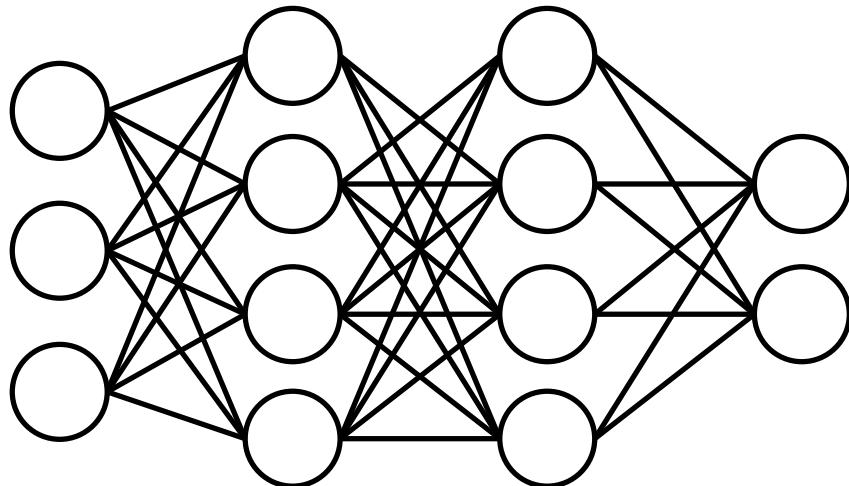
$$y = \mathbf{w}_L^T \sigma(\dots \mathbf{w}_2^T \sigma(\mathbf{w}_1^T \mathbf{x}))$$



small weights = sigmoid is in  
the linear range

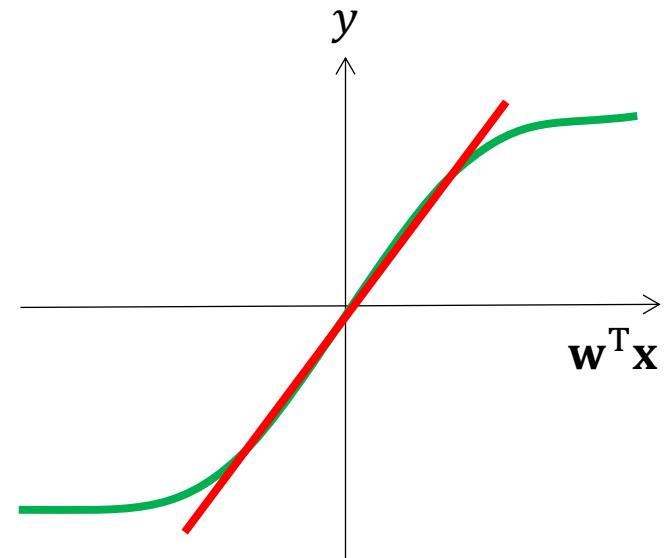
# Intuition

Why small weights reduce capacity?



$$y = \mathbf{w}_L^T (\dots \mathbf{w}_2^T (\mathbf{w}_1^T \mathbf{x})) = \mathbf{w}^T \mathbf{x}$$

behaves like single layer neural network



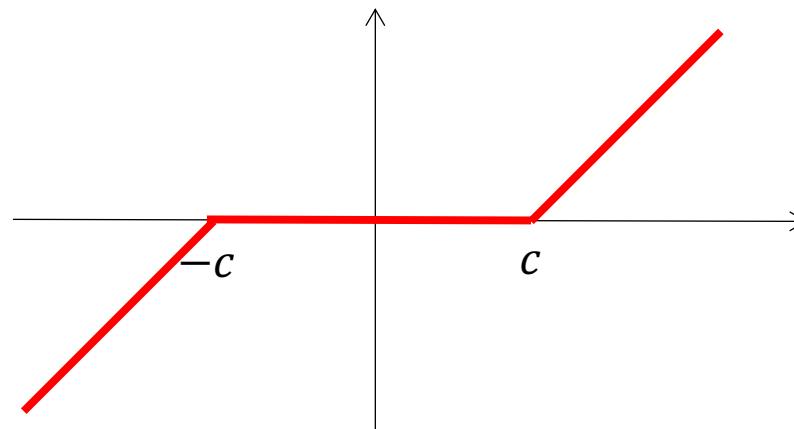
small weights = sigmoid is in  
the linear range

# $L_1$ Regularization

$$\hat{L}_R(\boldsymbol{\theta}) = \hat{L}(\boldsymbol{\theta}) + \beta \|\boldsymbol{\theta}\|_1$$

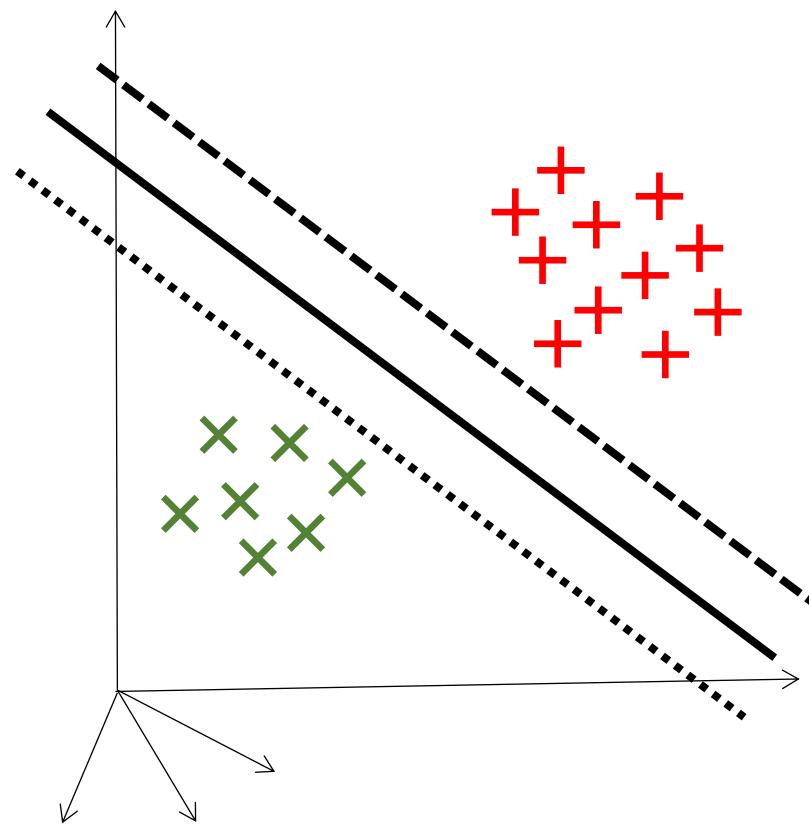
- Effect on solution: **shrinkage** (applied coordinate-wise) promoting sparsity

$$\boldsymbol{\theta}_R^* \approx \text{shrink}(\boldsymbol{\theta}^*) = \beta \text{sign}(\boldsymbol{\theta}^*) \max\{|\boldsymbol{\theta}^*| - c, 0\}$$

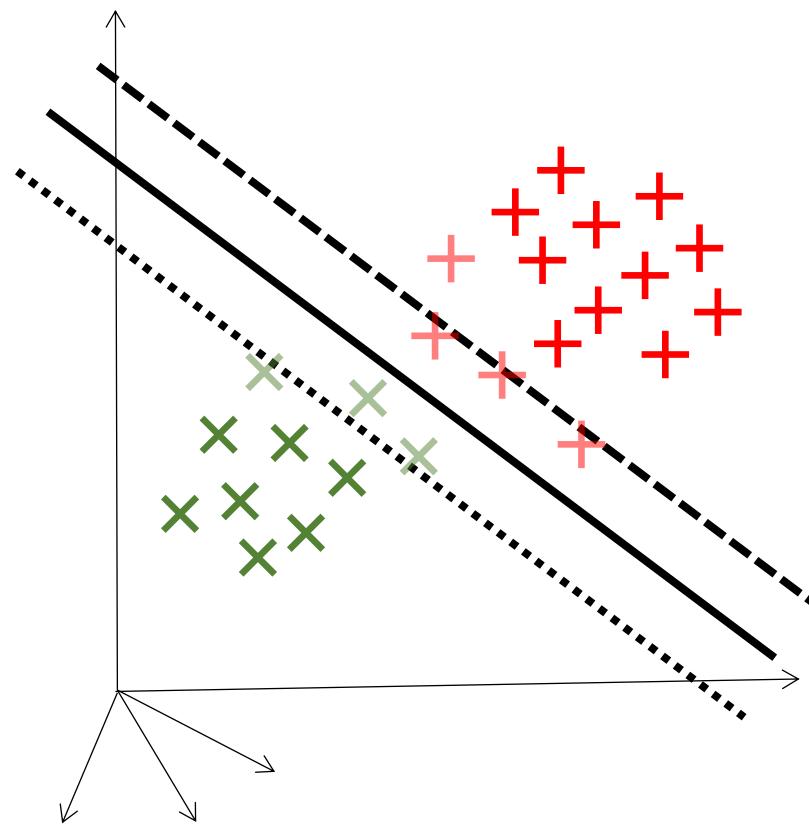


- Bayesian interpretation: Laplacian prior  $p(\boldsymbol{\theta}) \propto e^{\beta \|\boldsymbol{\theta}\|_1}$

# Noisy input



# Noisy input



Add noise to data

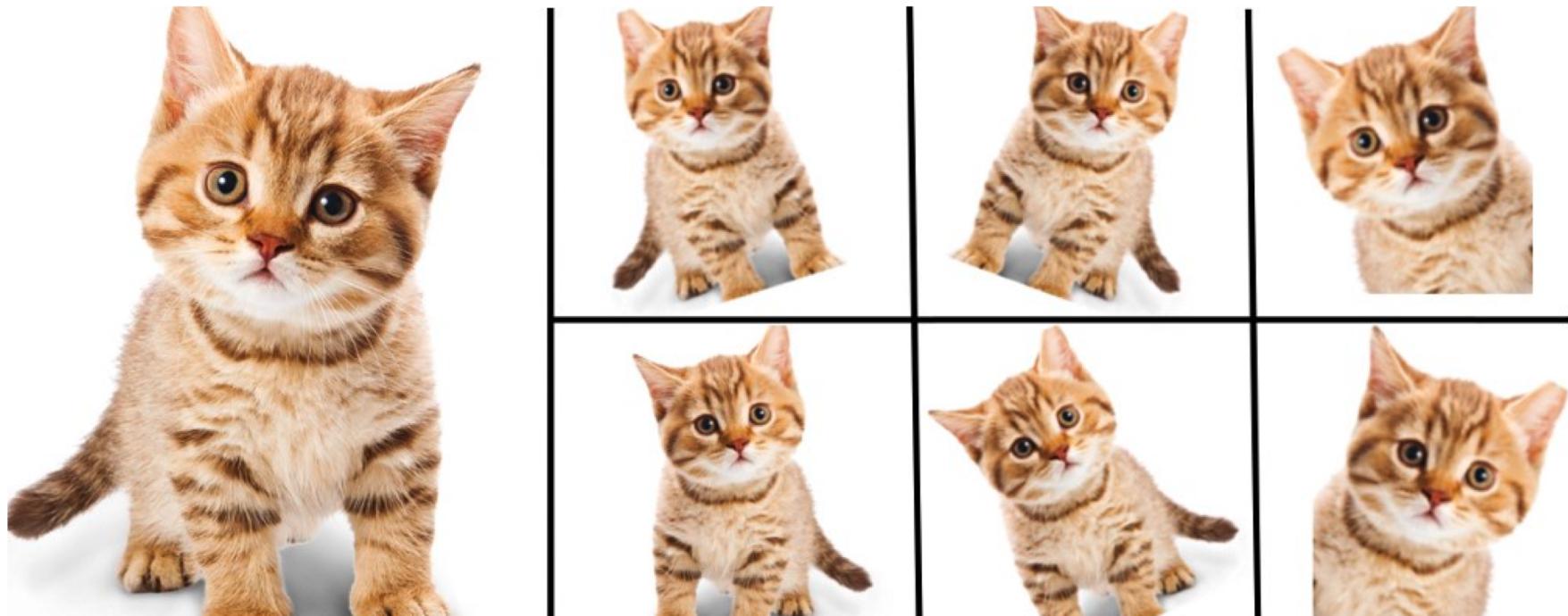
# Noisy input

- Add i.i.d. noise  $\epsilon \sim \mathcal{N}(0, \beta \mathbf{I})$  to input
- Loss becomes

$$\begin{aligned} L(\mathbf{w}) &= \mathbb{E}(\mathbf{w}^T(\mathbf{x} + \boldsymbol{\epsilon}) - y)^2 \\ &= \mathbb{E}(\mathbf{w}^T\mathbf{x} - y)^2 + 2\mathbb{E}\mathbf{w}^T\boldsymbol{\epsilon}(\mathbf{w}^T\mathbf{x} - y) + \mathbb{E}(\mathbf{w}^T\boldsymbol{\epsilon})^2 \\ &= \mathbb{E}(\mathbf{w}^T\mathbf{x} - y)^2 + \mathbf{w}^T\mathbb{E}(\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T)\mathbf{w} \\ &= \mathbb{E}(\mathbf{w}^T\mathbf{x} - y)^2 + \beta\|\mathbf{w}\|^2 \end{aligned}$$

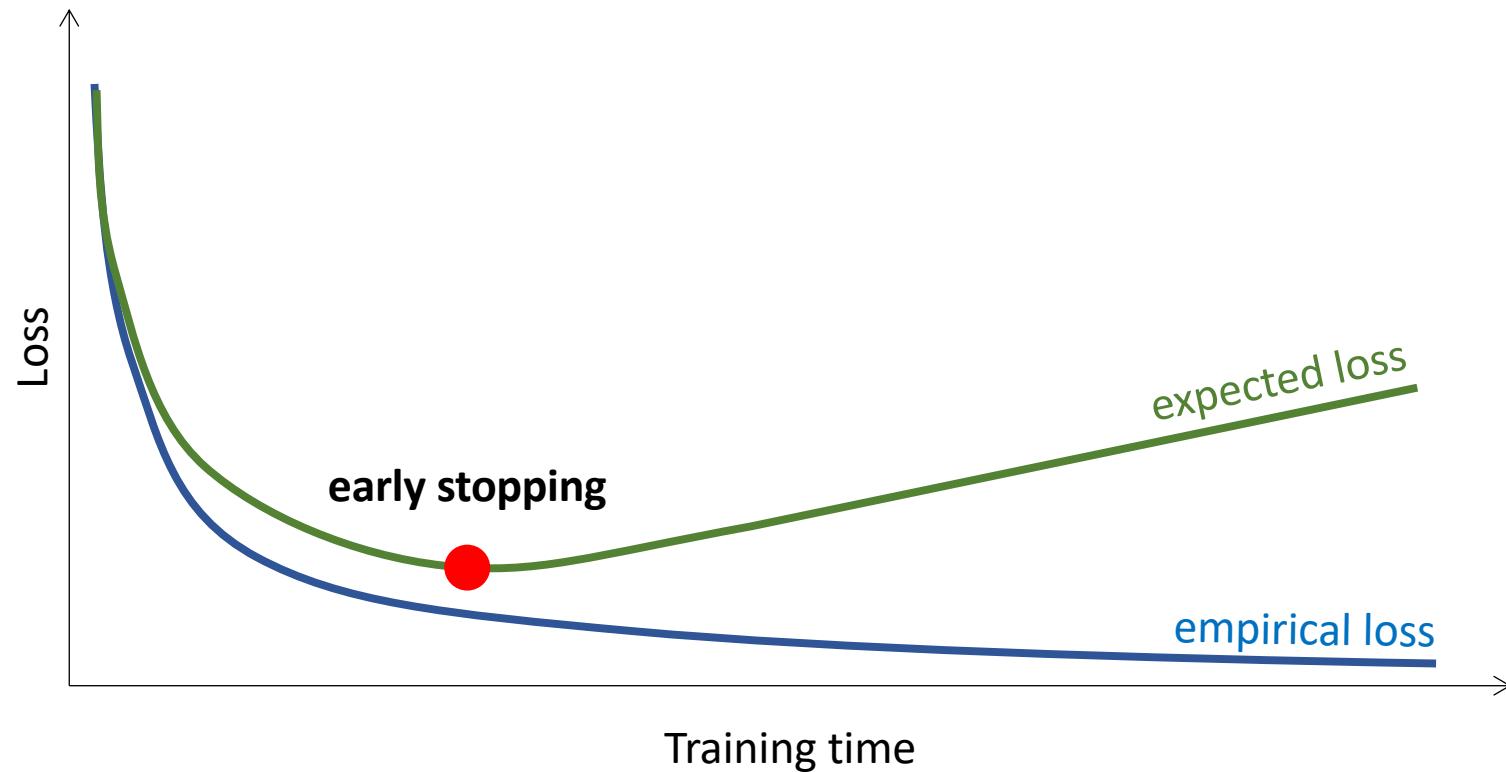
- Equivalent to weight decay!

# Data augmentation



Very common in computer vision

# Early stopping



# Early stopping

- Terminate before reaching convergence
- Effect similar to weight decay regularization (eigenvalue filtering)

# Early stopping as regularization

- Gradient of the loss around original problem optimum  $\theta^*$

$$\nabla_{\theta} \hat{L}(\theta) \approx \nabla_{\theta}^2 \hat{L}(\theta^*) (\theta - \theta^*)$$

- For small gradient descent step  $\alpha$ :

$$\begin{aligned}\theta^{(k+1)} &= \theta^{(k)} - \alpha \nabla_{\theta} \hat{L}(\theta^{(k)}) \\ &\approx \theta^{(k)} - \alpha \nabla_{\theta}^2 \hat{L}(\theta^*) (\theta^{(k)} - \theta^*)\end{aligned}$$

# Early stopping as regularization

- Gradient of the loss around original problem optimum  $\theta^*$

$$\nabla_{\theta} \hat{L}(\theta) \approx \nabla_{\theta}^2 \hat{L}(\theta^*) (\theta - \theta^*)$$

- For small gradient descent step  $\alpha$ :

$$\theta^{(k+1)} - \theta^* \approx \left( \mathbf{I} - \alpha \nabla_{\theta}^2 \hat{L}(\theta^*) \right) (\theta^{(k)} - \theta^*)$$

- Applying recursively from  $\theta^{(k)} = \mathbf{0}$ :

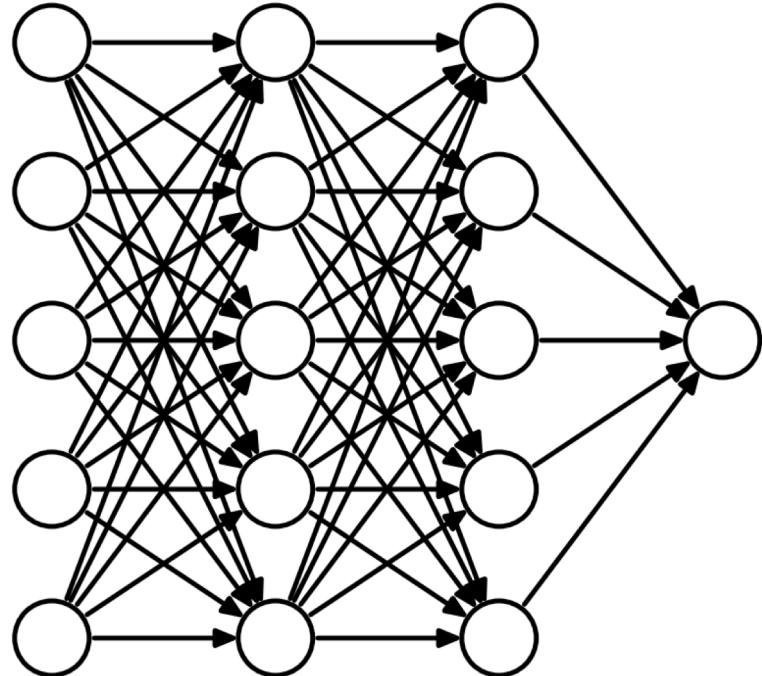
$$\begin{aligned}\theta^{(k+1)} &\approx \left( \mathbf{I} - \left( \mathbf{I} - \alpha \nabla_{\theta}^2 \hat{L}(\theta^*) \right)^{k+1} \right) \theta^* \\ &= \tau \left( \nabla_{\theta}^2 \hat{L}(\theta^*) \right) \theta^*\end{aligned}$$

- **Filter of  $\theta^*$**
- Bounds on stopping time

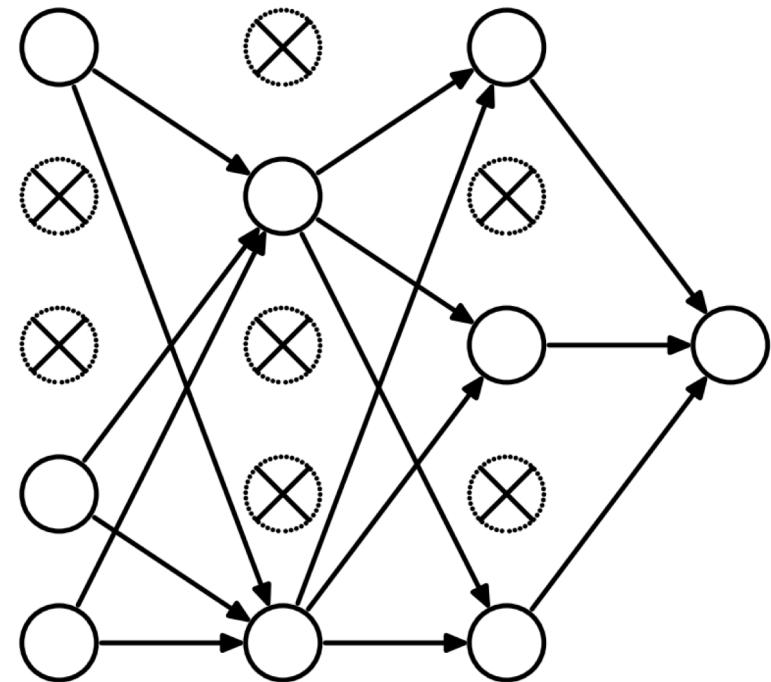
# Dropout



# Dropout



Without dropout

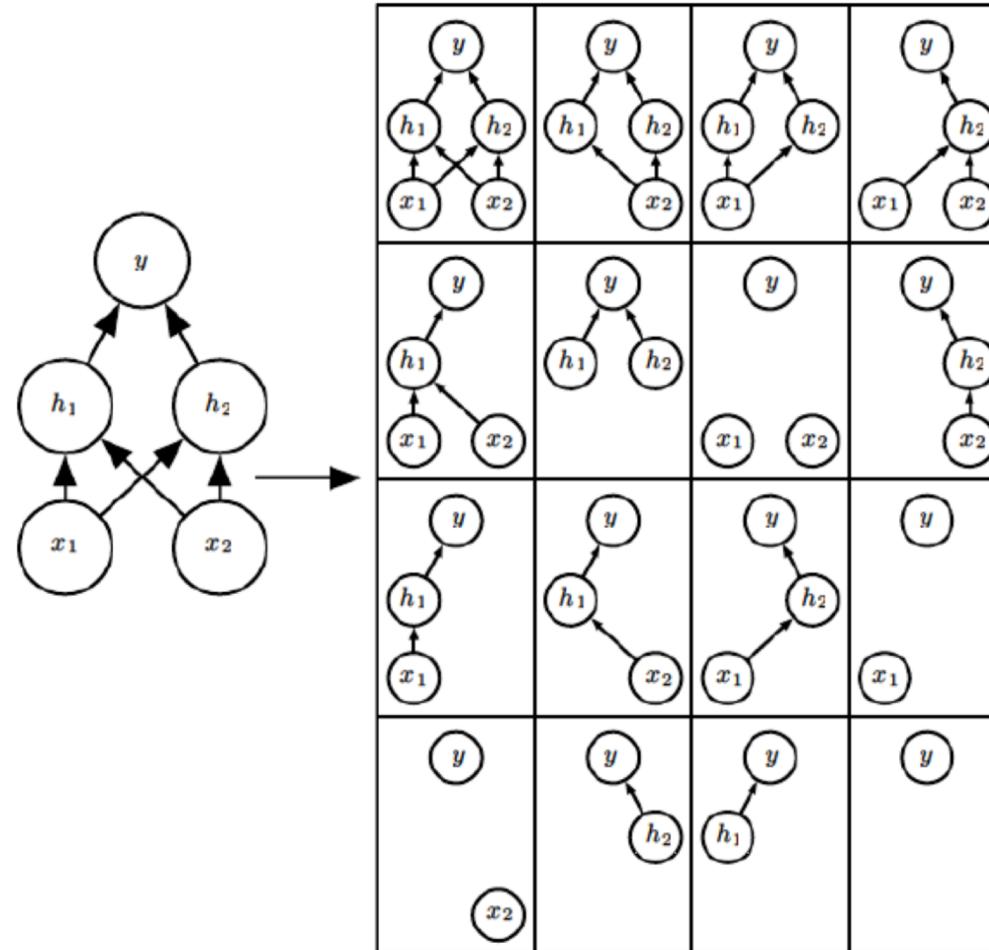


With dropout

# Dropout

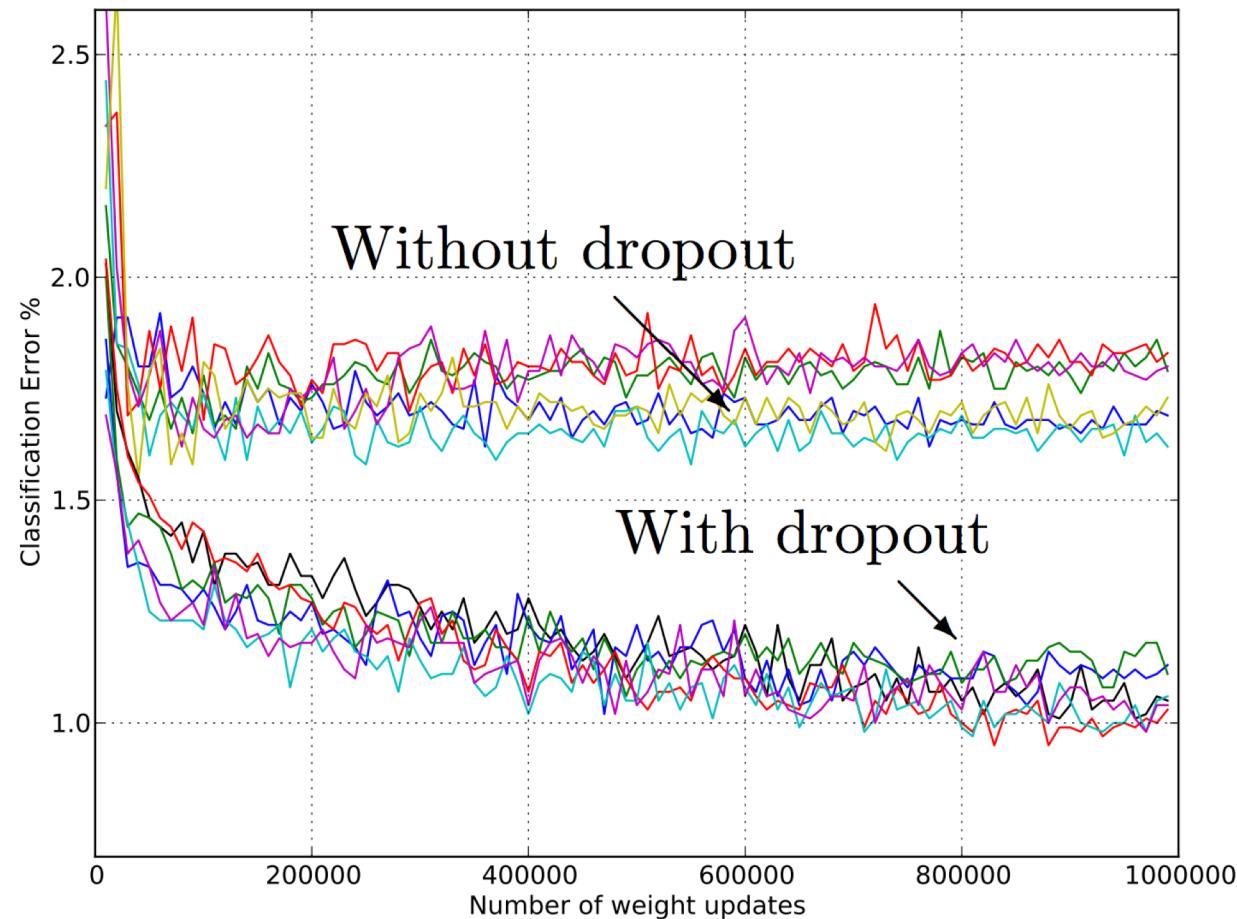
- During training, multiply output  $y_i$  of neuron  $i$  by a mask  $m_i \sim \text{Bernoulli}(p)$
- Explore various models with different connectivity
- “Ensemble learning”

# Dropout

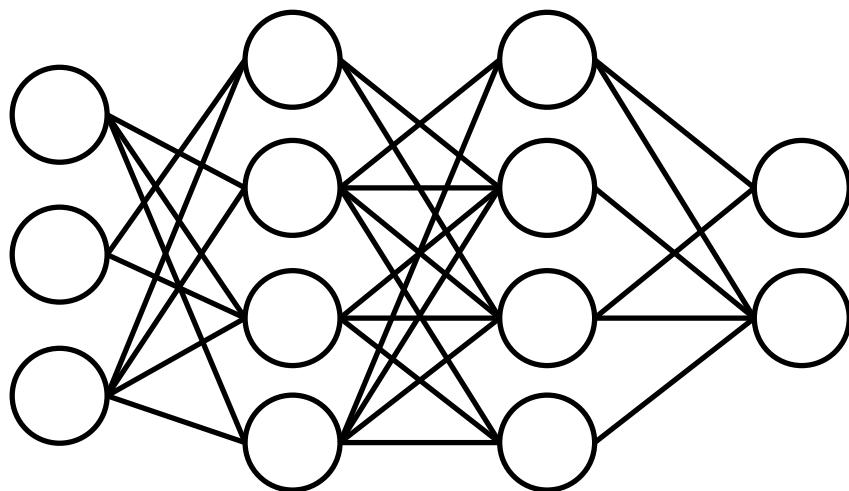


“Average on ensemble of models”

# Dropout

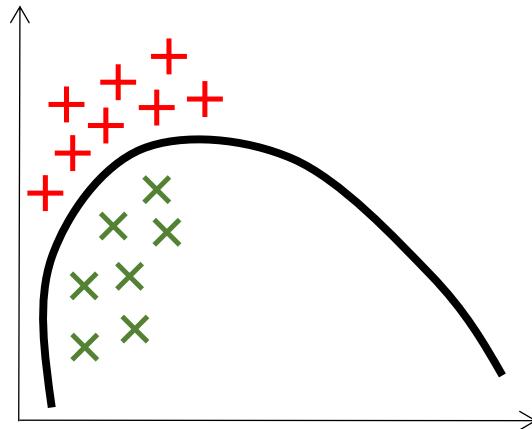


# DropConnect

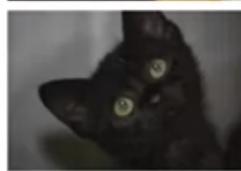


Turn off weights not neurons

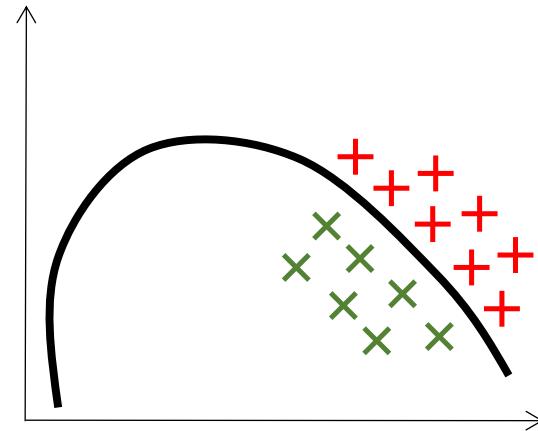
# Covariate shift



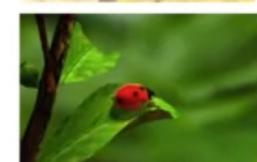
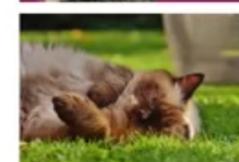
Training



Covariate  
shift

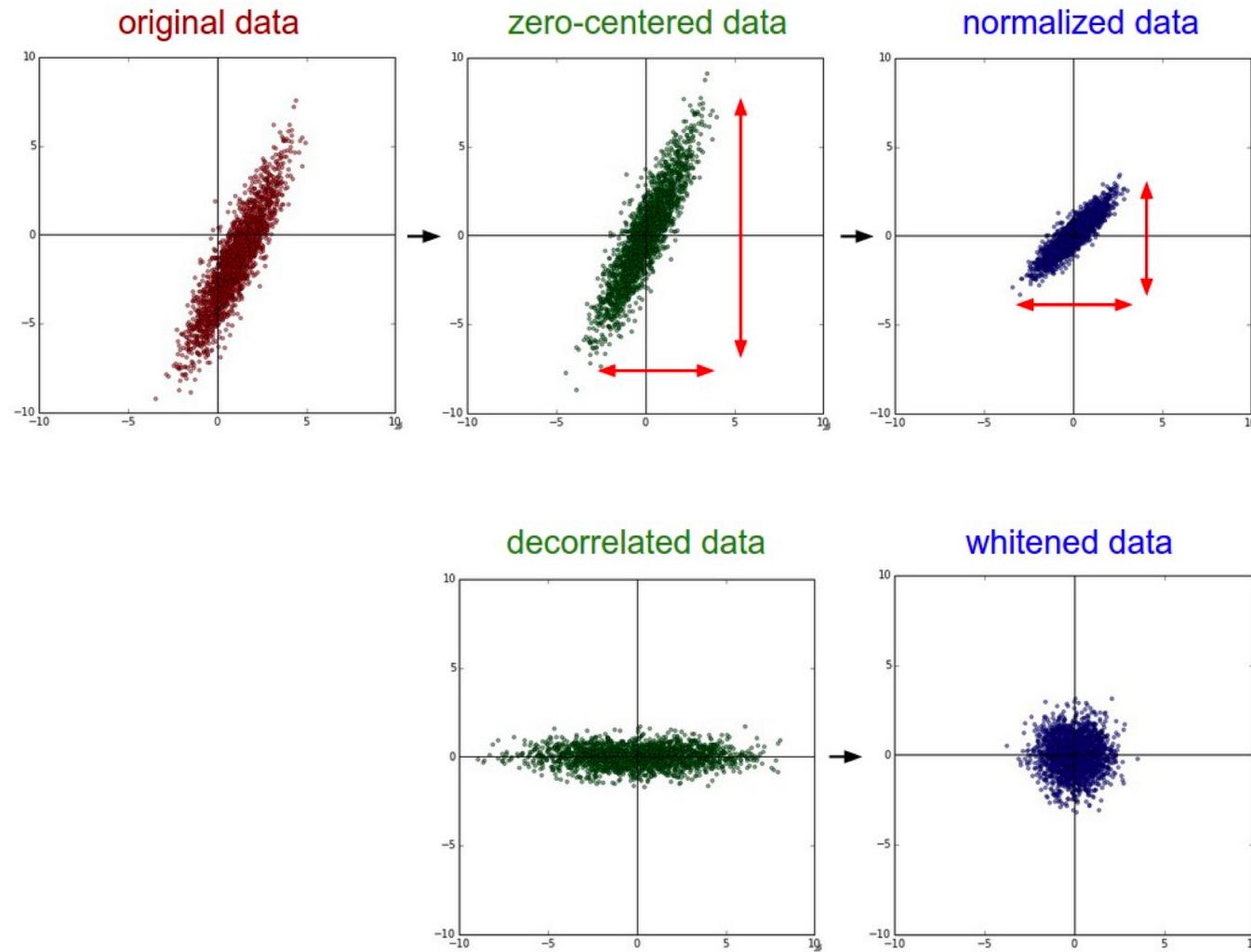


Testing



Example: A. Ng

# Input pre-processing



# Batch normalization

- Normalize each input dimension of a layer

$$\bar{x}_i = \frac{x_i - \mathbb{E}x_i}{\sqrt{\text{Var}x_i}}$$

where mean and variance are computed on a mini-batch

- Shift and scale normalized value by **learnable** parameters

$$y_i = \gamma_i \bar{x}_i + \beta_i$$