# Learning Dynamics of Linear Denoising Autoencoders

Arnu Pretorius, Steve Kroon and Herman Kamper
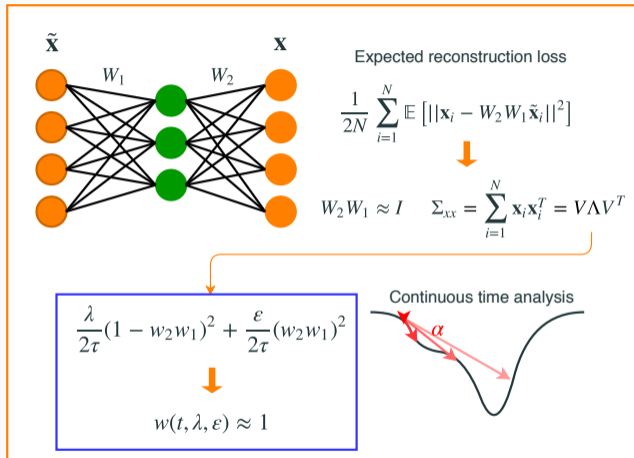
Stellenbosch University, South Africa

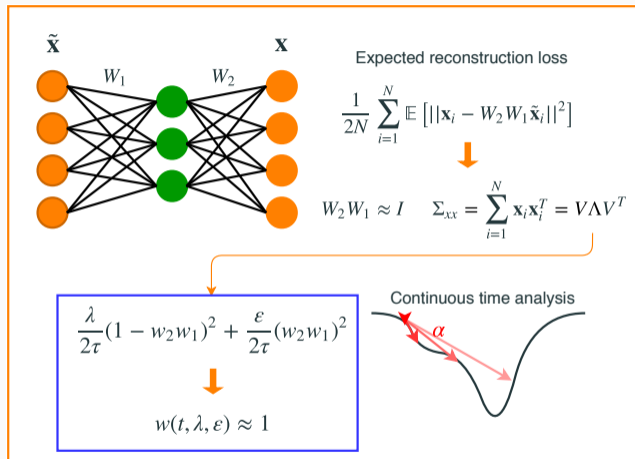$\tilde{\mathbf{x}}$ $W_1$ $W_2$ $\mathbf{x}$

Expected reconstruction loss

$$\frac{1}{2N} \sum_{i=1}^{N} \mathbb{E}\left[||\mathbf{x}_i - W_2 W_1 \tilde{\mathbf{x}}_i||^2\right]$$

$$W_2 W_1 \approx I \qquad \Sigma_{xx} = \sum_{i=1}^{N} \mathbf{x}_i \mathbf{x}_i^T = V\Lambda V^T$$

$$\frac{\lambda}{2\tau}(1 - w_2 w_1)^2 + \frac{\varepsilon}{2\tau}(w_2 w_1)^2$$

$$w(t, \lambda, \varepsilon) \approx 1$$
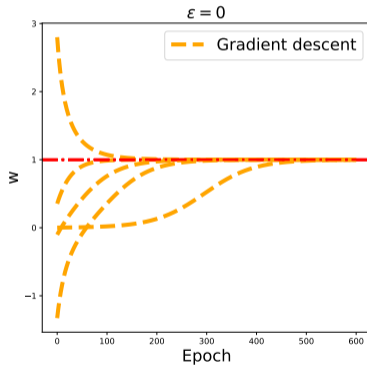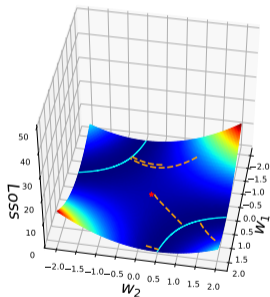
Continuous time analysis

$\alpha$

# Linear denoising autoencoders (DAE)



- *Exact solutions to the nonlinear dynamics of learning in deep linear neural networks*, Saxe, McClelland, Ganguli. ICLR, 2014.
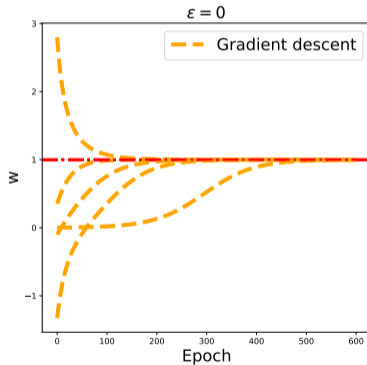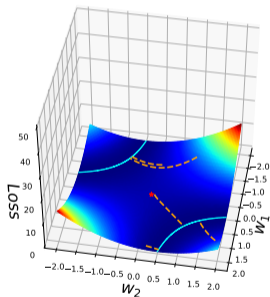
$$\ell_\varepsilon = \frac{\lambda}{2\tau}(1 - w_2 w_1)^2 + \frac{\varepsilon}{2\tau}(w_2 w_1)^2 \xrightarrow{\hspace{2cm}} \text{GD learning dynamics}$$
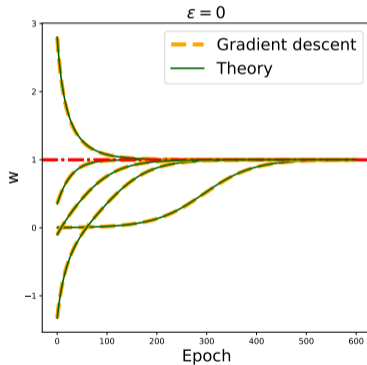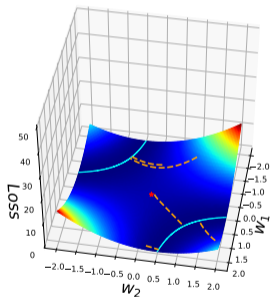
$$\ell_\varepsilon = \frac{\lambda}{2\tau}(1 - w_2 w_1)^2 + \frac{\varepsilon}{2\tau}(w_2 w_1)^2 \quad\longrightarrow\quad w(t, \lambda, \varepsilon) \quad\longrightarrow\quad \text{GD learning dynamics}$$

$$\ell_\varepsilon = \frac{\lambda}{2\tau}(1 - w_2 w_1)^2 + \frac{\varepsilon}{2\tau}(w_2 w_1)^2 \longrightarrow w(t, \lambda, \varepsilon) \longrightarrow \text{GD learning dynamics}$$
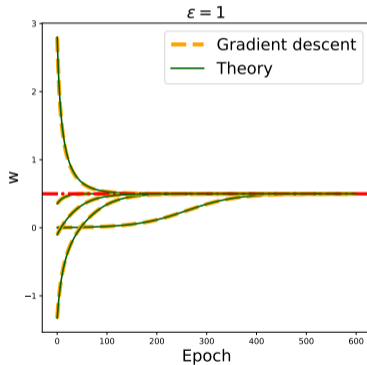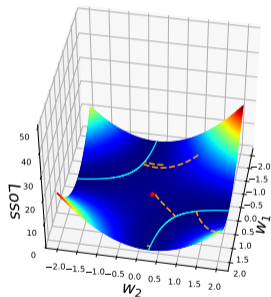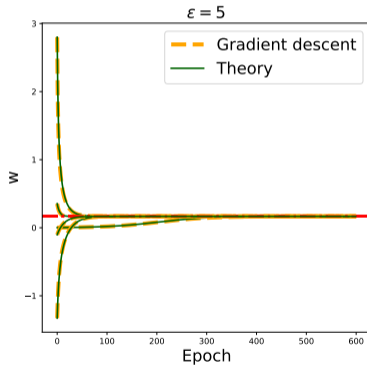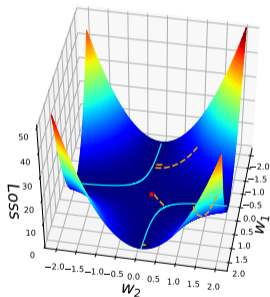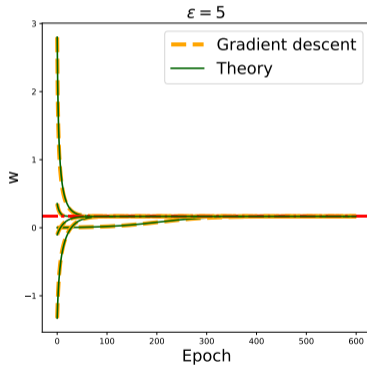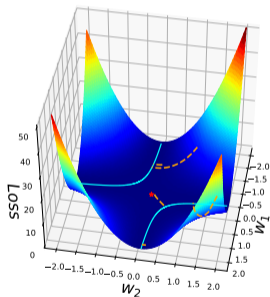
$$\ell_\varepsilon = \frac{\lambda}{2\tau}(1 - w_2 w_1)^2 + \frac{\varepsilon}{2\tau}(w_2 w_1)^2 \longrightarrow w(t, \lambda, \varepsilon) \longrightarrow \text{GD learning dynamics}$$

$$\ell_\varepsilon = \frac{\lambda}{2\tau}(1 - w_2 w_1)^2 + \frac{\varepsilon}{2\tau}(w_2 w_1)^2 \longrightarrow w(t, \lambda, \varepsilon) \longrightarrow \text{GD learning dynamics}$$
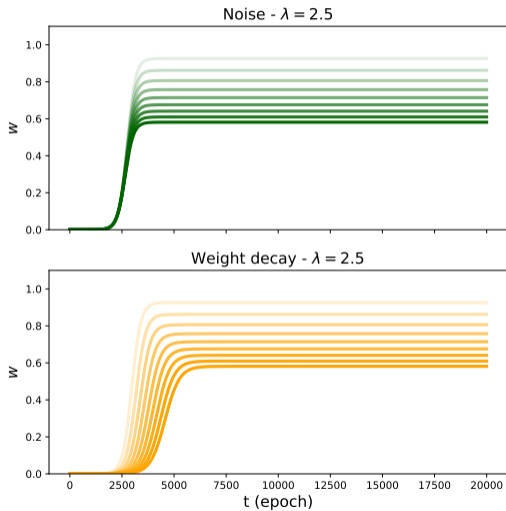
$$\ell_\varepsilon = \frac{\lambda}{2\tau}(1 - w_2 w_1)^2 + \frac{\varepsilon}{2\tau}(w_2 w_1)^2 \quad\longrightarrow\quad w(t, \lambda, \varepsilon) \quad\longrightarrow\quad \text{GD learning dynamics}$$
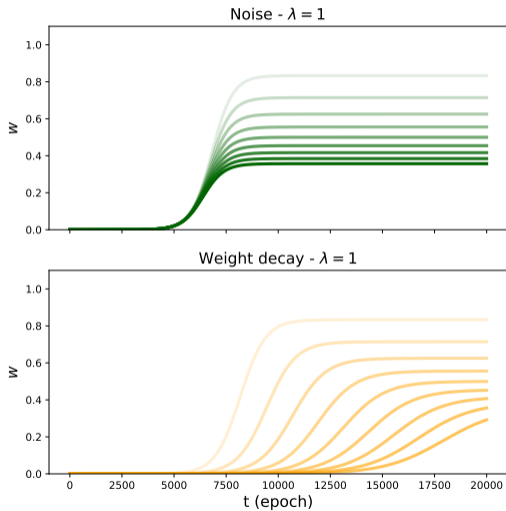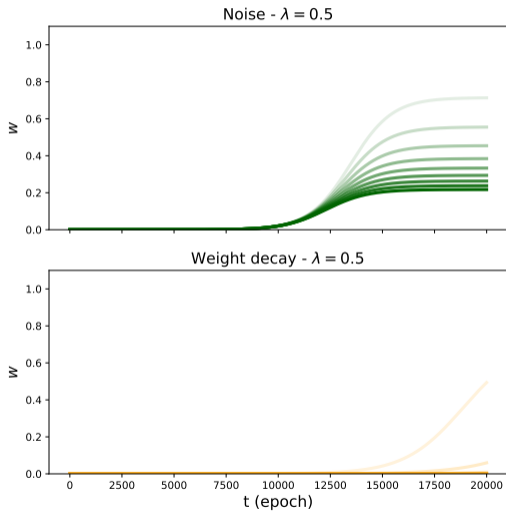


- Fixed point: $w^* = \frac{\lambda}{\lambda + \varepsilon}$
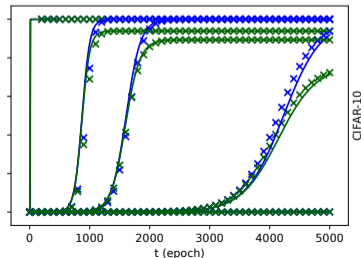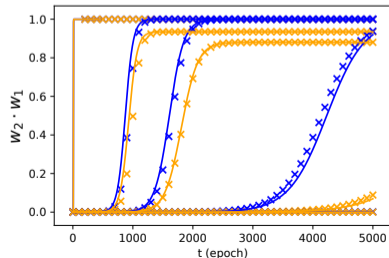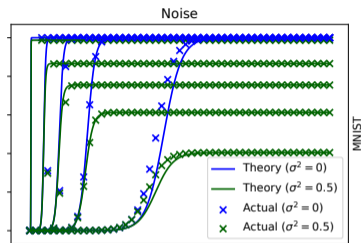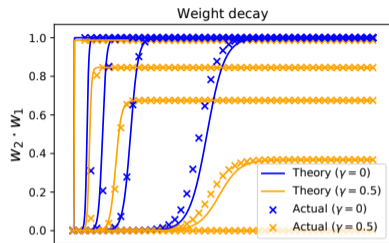
# The relationship between noise and weight decay

**Thank you for listening!**

**Source code to reproduce all the results**
https:
//github.com/arnupretorius/lindaedynamics_icml2018

## Optimal discrete time learning rates

- Ratio for DAE to WDAE:

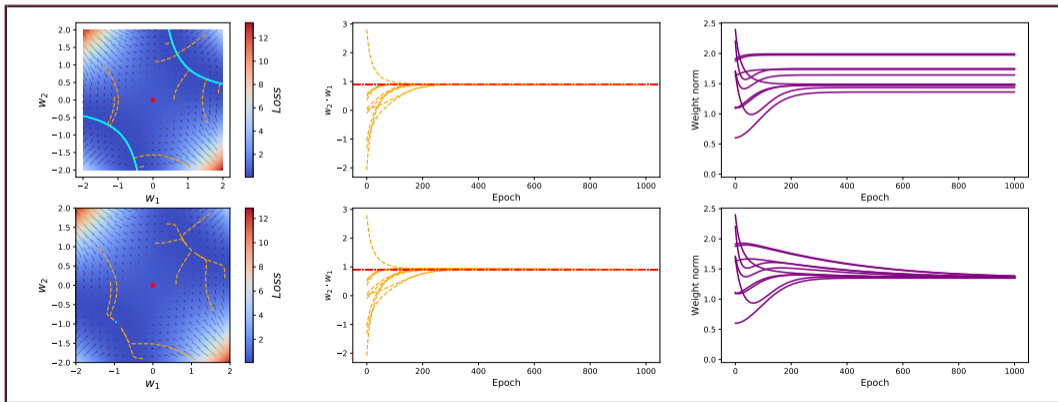$$R = \frac{2\lambda + \gamma}{2\lambda + 3\varepsilon}.$$

# The relationship between noise and weight decay
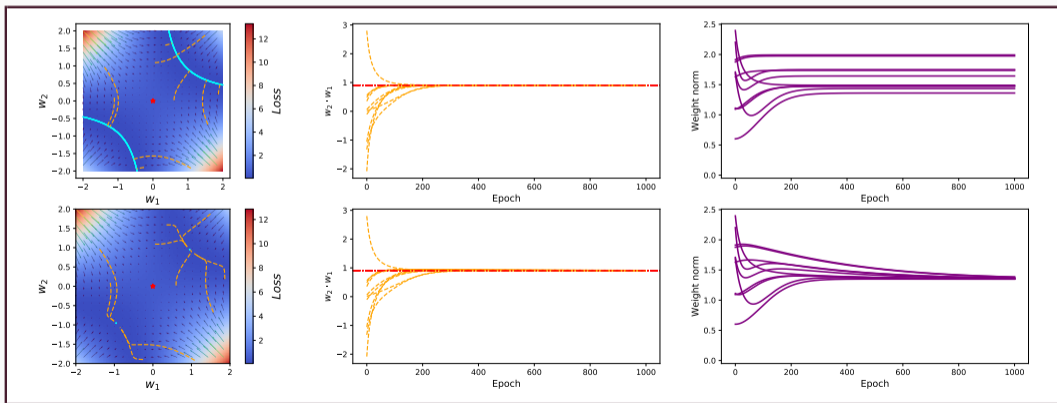
## Motivation for weight decay

- Smaller weights $\implies$ smoother models $\implies$ better generalisation

# The relationship between noise and weight decay

**Motivation for weight decay**

- Smaller weights $\implies$ smoother models $\implies$ better generalisation
- Small weight initialisation?

**Motivation for weight decay**

- Smaller weights $\implies$ smoother models $\implies$ better generalisation

- Small weight initialisation?