

## Weight shrinkage.

We do Ridge regression:

$$D\mathcal{L} = \frac{1}{N} \sum_n (\vec{w} \vec{x}_n - t_n) x_n + \lambda \vec{w}^2$$

So the basic GD step is:

$$\theta \mapsto \theta - \eta D\mathcal{L}$$

$$\text{or } \vec{w} \mapsto \vec{w} - \eta \lambda \vec{w} + -\eta \frac{1}{N} \sum_{n=1}^N (w x_n - t_n) x_n$$

$$\vec{w} (1 - \eta \lambda) - \underbrace{\quad}_{\text{normal MSE term.}}$$

geometrical decay  
with a rate  $1 - \eta \lambda < 1$ .

Hence the name "weight decay"

or "parameter shrinkage"

(a geometric series decays exponentially fast,  
so if the  $-\frac{\eta}{N} \sum_n (\quad) x_n$  term does not  
push  $\vec{w}$ , it goes towards 0).