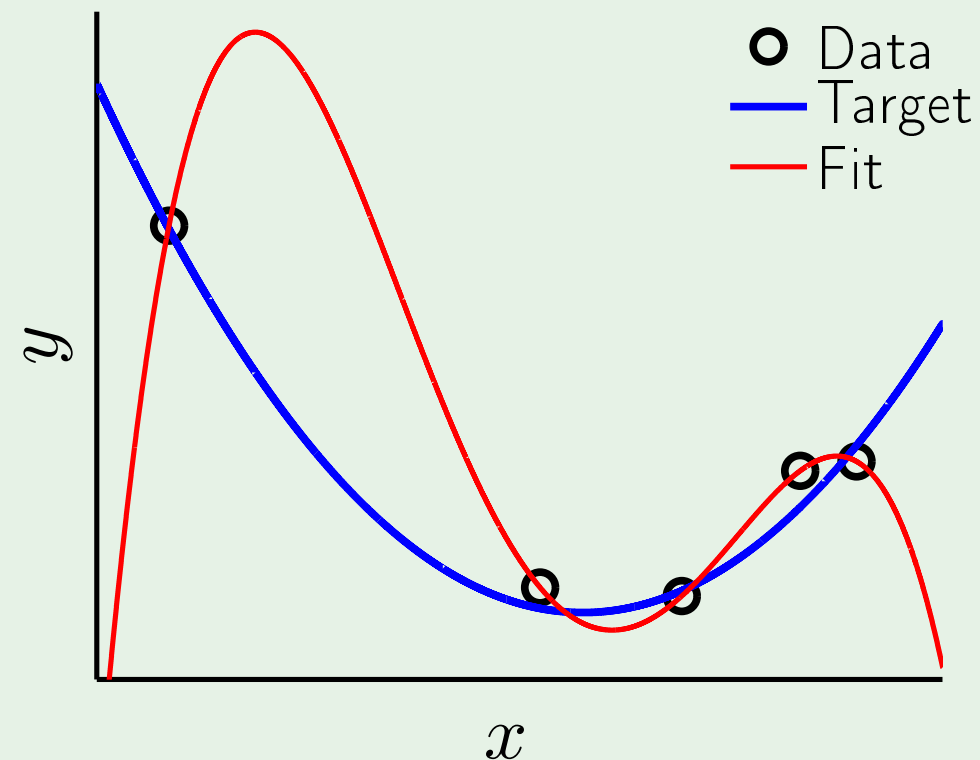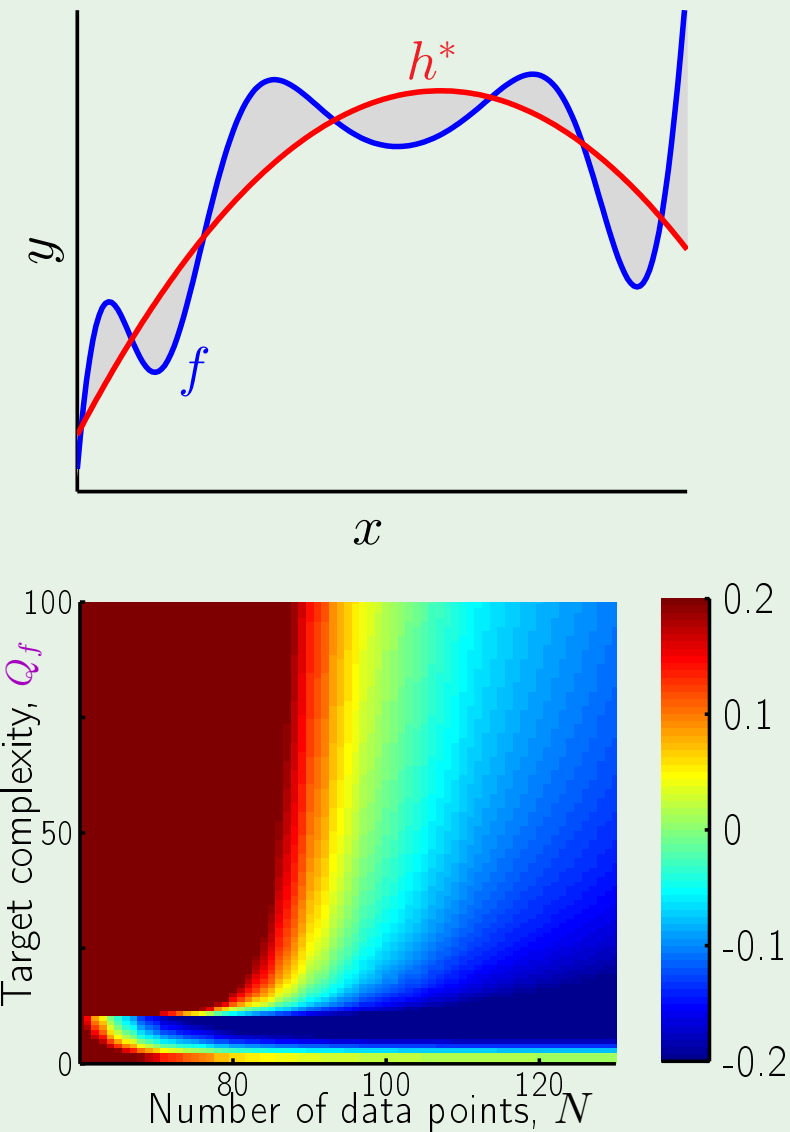# Review of Lecture 11

- **Overfitting**

  Fitting the data more than is warranted



  VC allows it; doesn't predict it

Fitting the noise, stochastic/deterministic

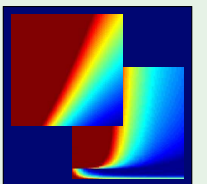- **Deterministic noise**

# Learning From Data

Yaser S. Abu-Mostafa
*California Institute of Technology*

Lecture 12: **Regularization**

# Outline

- Regularization - informal

- Regularization - formal

- Weight decay

- Choosing a regularizer

# Two approaches to regularization

## Mathematical:

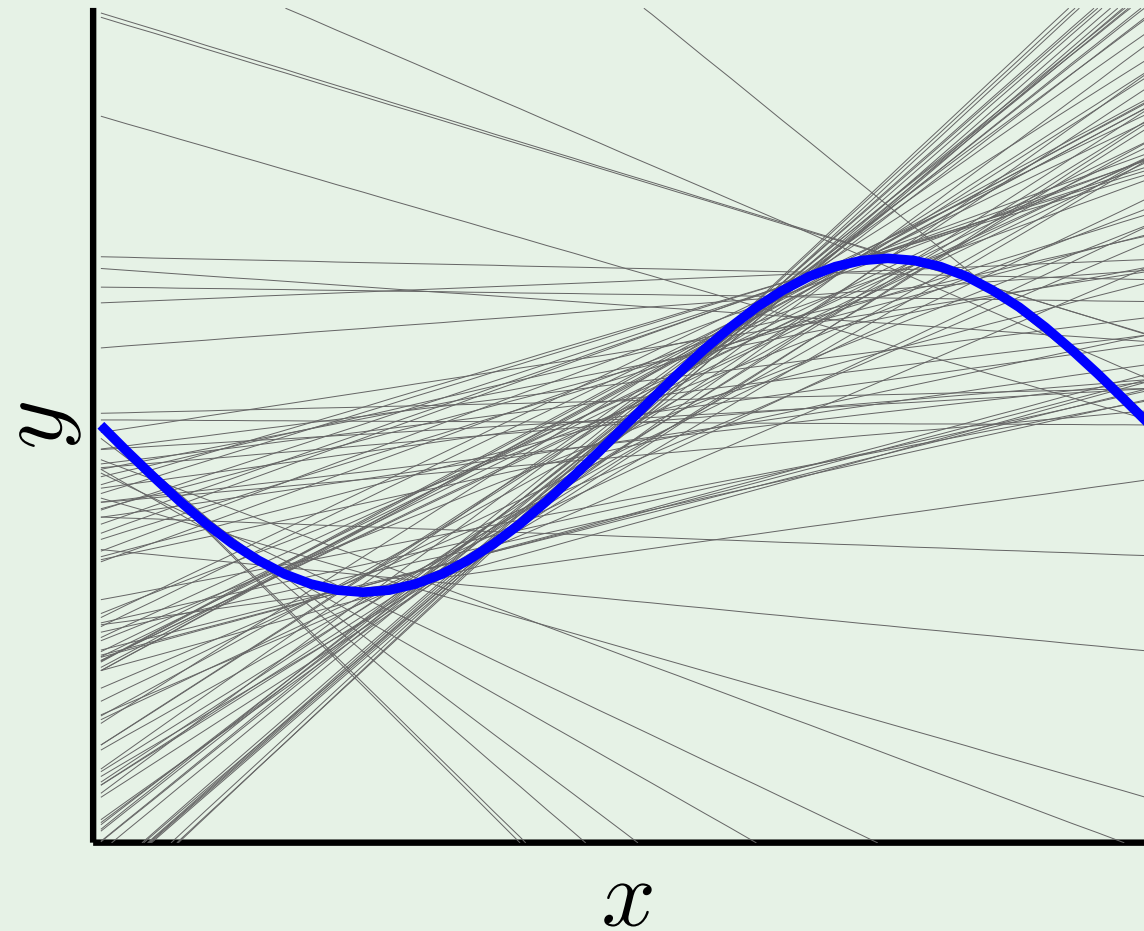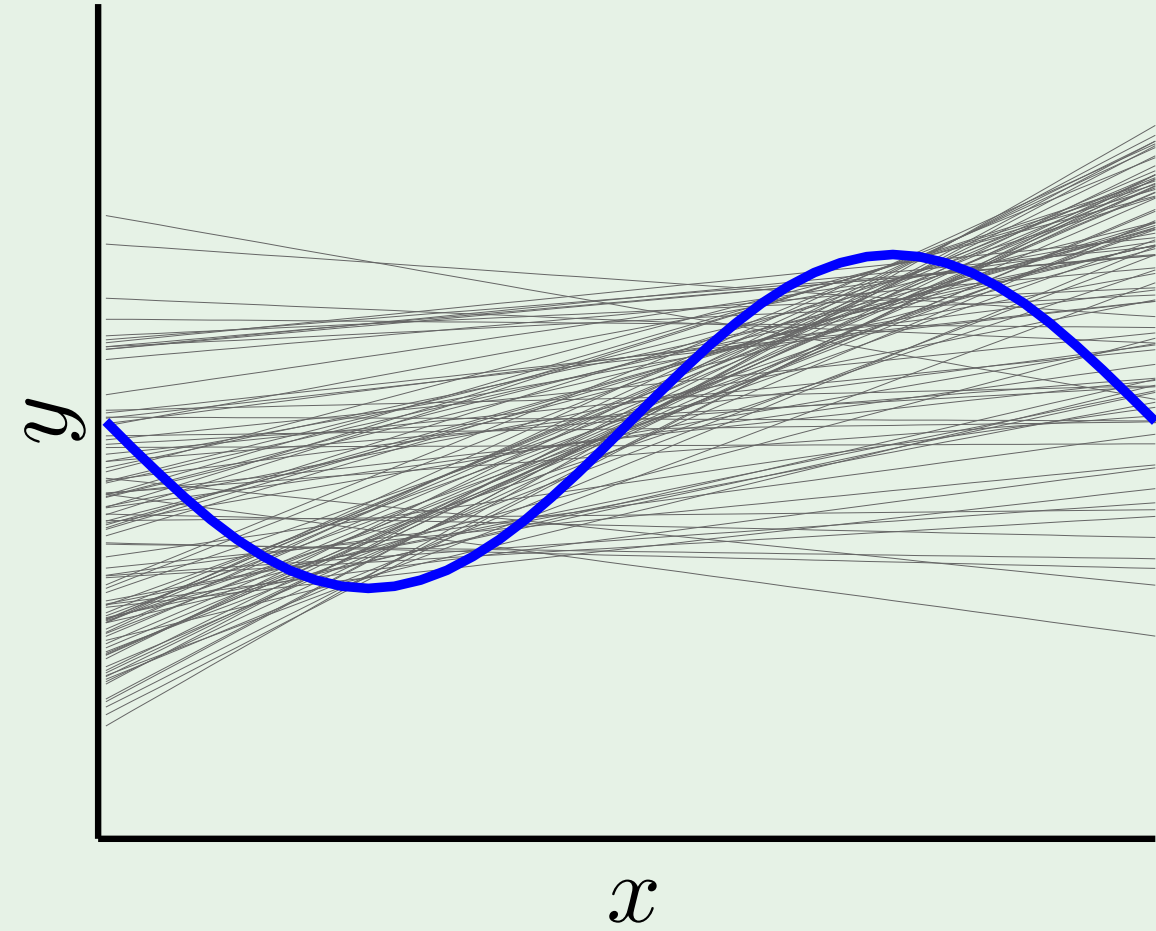Ill-posed problems in function approximation

## Heuristic:

Handicapping the minimization of $E_{\text{in}}$

# A familiar example



without regularization

with regularization

# and the winner is ...



without regularization

$\bar{g}(x)$

$\sin(\pi x)$

bias $= \mathbf{0.21}$    var $= \mathbf{1.69}$

with regularization

$\bar{g}(x)$

$\sin(\pi x)$

bias $= \mathbf{0.23}$    var $= \mathbf{0.33}$

# The polynomial model

$\mathcal{H}_{\mathrm{Q}}$: polynomials of order $Q$        linear regression in $\mathcal{Z}$ space

$$\mathbf{z} = \begin{bmatrix} 1 \\ L_1(x) \\ \vdots \\ L_{\mathrm{Q}}(x) \end{bmatrix} \qquad \mathcal{H}_{\mathrm{Q}} = \left\{ \sum_{q=0}^{Q} w_q \, L_q(x) \right\}$$

Legendre polynomials:



| $L_1$ | $L_2$ | $L_3$ | $L_4$ | $L_5$ |
|---|---|---|---|---|
| $x$ | $\frac{1}{2}(3x^2 - 1)$ | $\frac{1}{2}(5x^3 - 3x)$ | $\frac{1}{8}(35x^4 - 30x^2 + 3)$ | $\frac{1}{8}(63x^5 \cdots)$ |

# Unconstrained solution

Given $(x_1, y_1), \cdots, (x_N, y_n)$ $\longrightarrow$ $(\mathbf{z}_1, y_1), \cdots, (\mathbf{z}_N, y_n)$

Minimize $E_{\text{in}}(\mathbf{w}) = \dfrac{1}{N} \displaystyle\sum_{n=1}^{N} (\mathbf{w}^{\intercal}\mathbf{z}_n - y_n)^2$

Minimize $\dfrac{1}{N} (\mathrm{Z}\mathbf{w} - \mathbf{y})^{\intercal}(\mathrm{Z}\mathbf{w} - \mathbf{y})$

$$\mathbf{w}_{\text{lin}} = (\mathrm{Z}^{\intercal}\mathrm{Z})^{-1}\mathrm{Z}^{\intercal}\mathbf{y}$$

pseudo inverse

# Constraining the weights

Hard constraint:      $\mathcal{H}_2$ is constrained version of $\mathcal{H}_{10}$      with $w_q = 0$ for $q > 2$

Softer version:      $\displaystyle\sum_{q=0}^{Q} w_q^2 \leq C$      "**soft-order**" constraint

Minimize   $\frac{1}{N} (\mathrm{Z}\mathbf{w} - \mathbf{y})^\top (\mathrm{Z}\mathbf{w} - \mathbf{y})$

subject to:   $\mathbf{w}^\top \mathbf{w} \leq C$

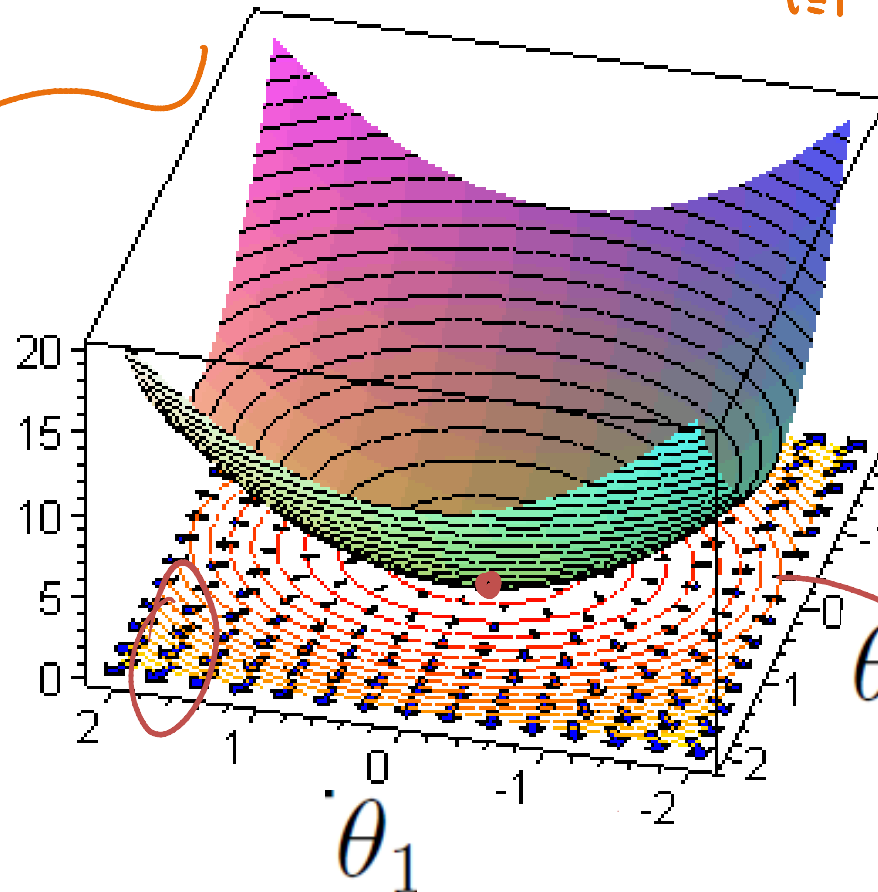Solution:  $\mathbf{w}_{\mathrm{reg}}$  instead of  $\mathbf{w}_{\mathrm{lin}}$

# Optimization approach

Our aim is to minimise the quadratic cost between the output labels and the model predictions

$$J(\boldsymbol{\theta}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) = \sum_{i=1}^{n} (y_i - \mathbf{x}_i^T \boldsymbol{\theta})^2 \; z \; \sum_{i=1}^{n} \left( Y_i - \hat{Y}_i \right)^2$$

$n \times d$

$h \times 1$  $d \times 1$  $|x|$

$|x|$

$f(\theta) = a\theta^2 + b\theta + c;$
quadratic graph
in 2D

$\hat{Y}_i = 1\theta_1 + X_i\theta_2$

$J(\theta_1, \theta_2)$

$\theta^2 <= C$
or $2\theta^2 <= 2C$
or $\theta^2 + \theta^2 <= sqrt(2C)^2$



$\theta_2$

$\theta_1$

contour lines

# Solving for $\mathbf{w}_{\mathrm{reg}}$

Minimize $\quad E_{\mathrm{in}}(\mathbf{w}) \;=\; \frac{1}{N}\,(\mathbf{Zw}-\mathbf{y})^{\top}(\mathbf{Zw}-\mathbf{y})$

$\qquad$ subject to: $\quad \mathbf{w}^{\top}\mathbf{w} \le C$

derive: W$^{\mathsf{T}}$W=> W²;
der(W²)=2W

$\nabla E_{\mathrm{in}}(\mathbf{w}_{\mathrm{reg}}) \;\propto\; -\mathbf{w}_{\mathrm{reg}}$

verticle cut

Math Theorem: gradients are always normal to level curves.

$\qquad\quad =\; -2\frac{\lambda}{N}\mathbf{w}_{\mathrm{reg}}$

$\nabla E_{\mathrm{in}}(\mathbf{w}_{\mathrm{reg}}) + 2\frac{\lambda}{N}\mathbf{w}_{\mathrm{reg}} \;=\; \mathbf{0}$

Minimize $\quad E_{\mathrm{in}}(\mathbf{w}) + \frac{\lambda}{N}\mathbf{w}^{\top}\mathbf{w}$

$\boxed{C \uparrow \quad \lambda \downarrow}$



$E_{\mathrm{in}} = \mathrm{const.}$

$\mathbf{w}_{\mathrm{lin}}$

normal

$\mathbf{w}$

$\nabla E_{\mathrm{in}}$

$\mathbf{w}^{\mathsf{T}}\mathbf{w} = C$

**Paraboloid Surface with Constraint Sphere**
$(\theta\,\theta \leq C)$

$E_n(\theta))$

$\theta_1$

$\theta_1$

$(\theta\,\mathsf{T}\,\theta\ C)$

# Augmented error

Minimizing $\quad E_{\text{aug}}(\mathbf{w}) = E_{\text{in}}(\mathbf{w}) + \frac{\lambda}{N}\mathbf{w}^{\mathsf{T}}\mathbf{w}$

$$= \frac{1}{N}\left(\mathbf{Z}\mathbf{w} - \mathbf{y}\right)^{\mathsf{T}}(\mathbf{Z}\mathbf{w} - \mathbf{y}) + \frac{\lambda}{N}\mathbf{w}^{\mathsf{T}}\mathbf{w} \qquad \text{unconditionally}$$

$-$ solves $-$

Minimizing $\quad E_{\text{in}}(\mathbf{w}) = \frac{1}{N}\left(\mathbf{Z}\mathbf{w} - \mathbf{y}\right)^{\mathsf{T}}(\mathbf{Z}\mathbf{w} - \mathbf{y})$

subject to: $\quad \mathbf{w}^{\mathsf{T}}\mathbf{w} \leq C \qquad\qquad \longleftarrow$ VC formulation

# The solution

Minimize $\qquad E_{\text{aug}}(\mathbf{w}) \;=\; E_{\text{in}}(\mathbf{w}) + \frac{\lambda}{N}\mathbf{w}^{\mathsf{T}}\mathbf{w}$

$$= \; \frac{1}{N}\left((\mathbf{Z}\mathbf{w}-\mathbf{y})^{\mathsf{T}}(\mathbf{Z}\mathbf{w}-\mathbf{y}) \;+\; \lambda\,\mathbf{w}^{\mathsf{T}}\mathbf{w}\right)$$

$\nabla E_{\text{aug}}(\mathbf{w}) \;=\; \mathbf{0} \qquad \Longrightarrow \qquad \mathbf{Z}^{\mathsf{T}}(\mathbf{Z}\mathbf{w}-\mathbf{y}) + \lambda\mathbf{w} = \mathbf{0}$

zzw-zy+kw
w(zz+kI) = zy
w = (zz+kI)⁻¹.zy

$$\boxed{\;\mathbf{w}_{\text{reg}} \;=\; (\mathbf{Z}^{\mathsf{T}}\mathbf{Z} + \lambda\mathbf{I})^{-1}\,\mathbf{Z}^{\mathsf{T}}\mathbf{y}\;}$$

(with regularization)

as opposed to $\qquad \mathbf{w}_{\text{lin}} \;=\; (\mathbf{Z}^{\mathsf{T}}\mathbf{Z})^{-1}\mathbf{Z}^{\mathsf{T}}\mathbf{y}$ $\qquad$ (without regularization)

pinv

# The result

Minimizing $\quad E_{\text{in}}(\mathbf{w}) + \dfrac{\lambda}{N}\,\mathbf{w}^{\mathsf{T}}\mathbf{w}\quad$ for different $\lambda$'s:



| $\lambda = 0$ | $\lambda = 0.0001$ | $\lambda = 0.01$ | $\lambda = 1$ |

$W_{\text{reg}} = W_{\text{lin}}$

$W_{\text{reg}}$ is not reaching $W_{\text{lin}}$

**overfitting** $\quad\longrightarrow\quad\longrightarrow\quad\longrightarrow\quad\longrightarrow\quad$ **underfitting**

# Weight 'decay'

Minimizing $\quad E_{\text{in}}(\mathbf{w}) + \frac{\lambda}{N} \mathbf{w}^{\mathsf{T}}\mathbf{w} \quad$ is called weight *decay*. Why?

Gradient descent:

$$\mathbf{w}(t+1) = \mathbf{w}(t) \ - \ \eta \, \nabla E_{\text{in}}\left(\mathbf{w}(t)\right) \ - \ 2\,\eta\,\frac{\lambda}{N}\,\mathbf{w}(t)$$

$$= \mathbf{w}(t)\,(1 - 2\eta\frac{\lambda}{N}) - \eta\,\nabla E_{\text{in}}\left(\mathbf{w}(t)\right)$$

Applies in neural networks:

$$\mathbf{w}^{\mathsf{T}}\mathbf{w} = \sum_{l=1}^{L} \sum_{i=0}^{d^{(l-1)}} \sum_{j=1}^{d^{(l)}} \left(w_{ij}^{(l)}\right)^{2}$$

# Variations of weight decay

Emphasis of certain weights:
$$\sum_{q=0}^{Q} \gamma_q \; w_q^2$$

Examples:

$$\gamma_q = 2^q \implies \text{low-order fit}$$

$$\gamma_q = 2^{-q} \implies \text{high-order fit}$$

Neural networks: different layers get different $\gamma$'s

**Tikhonov regularizer:** $\mathbf{w}^\mathsf{T} \Gamma^\mathsf{T} \Gamma \mathbf{w}$
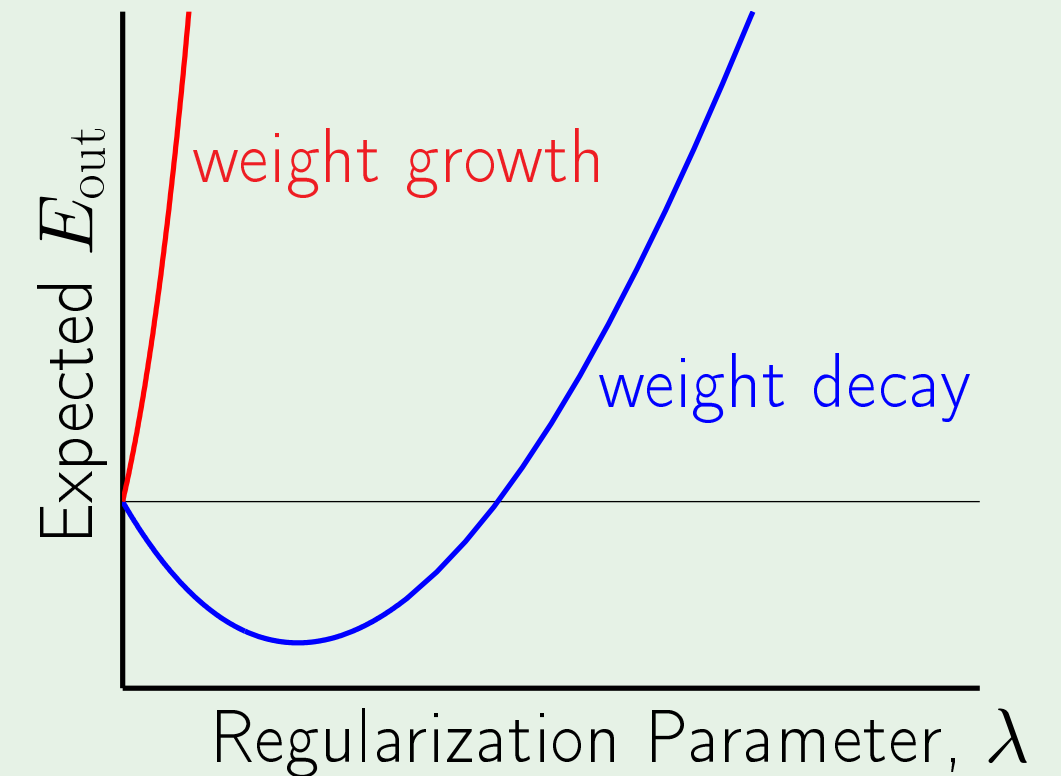
# Even weight growth!

We 'constrain' the weights to be large - bad!

## Practical rule:

    stochastic noise is 'high-frequency'

    deterministic noise is also non-smooth

$\Longrightarrow$   constrain learning towards smoother hypotheses

# General form of augmented error

Calling the regularizer $\Omega = \Omega(h)$, we minimize

$$E_{\text{aug}}(h) = E_{\text{in}}(h) + \frac{\lambda}{N}\Omega(h)$$

Rings a bell? $\downarrow\downarrow$

$$E_{\text{out}}(h) \leq E_{\text{in}}(h) + \Omega(\mathcal{H})$$

$E_{\text{aug}}$ is better than $E_{\text{in}}$ as a proxy for $E_{\text{out}}$

# Outline

- Regularization - informal

- Regularization - formal

- Weight decay

- <span style="color:blue">Choosing a regularizer</span>

# The perfect regularizer $\Omega$

Constraint in the 'direction' of the target function  (going in circles ☺)

Guiding principle:

Direction of **smoother** or "simpler"

Chose a bad $\Omega$?
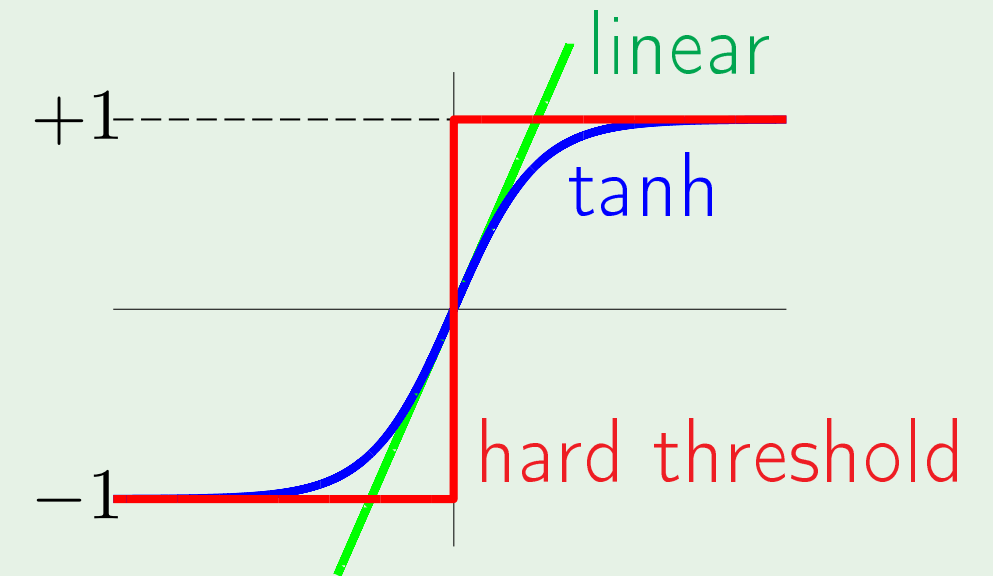
We still have $\lambda$!

# Neural-network regularizers

**Weight decay:** From linear to logical

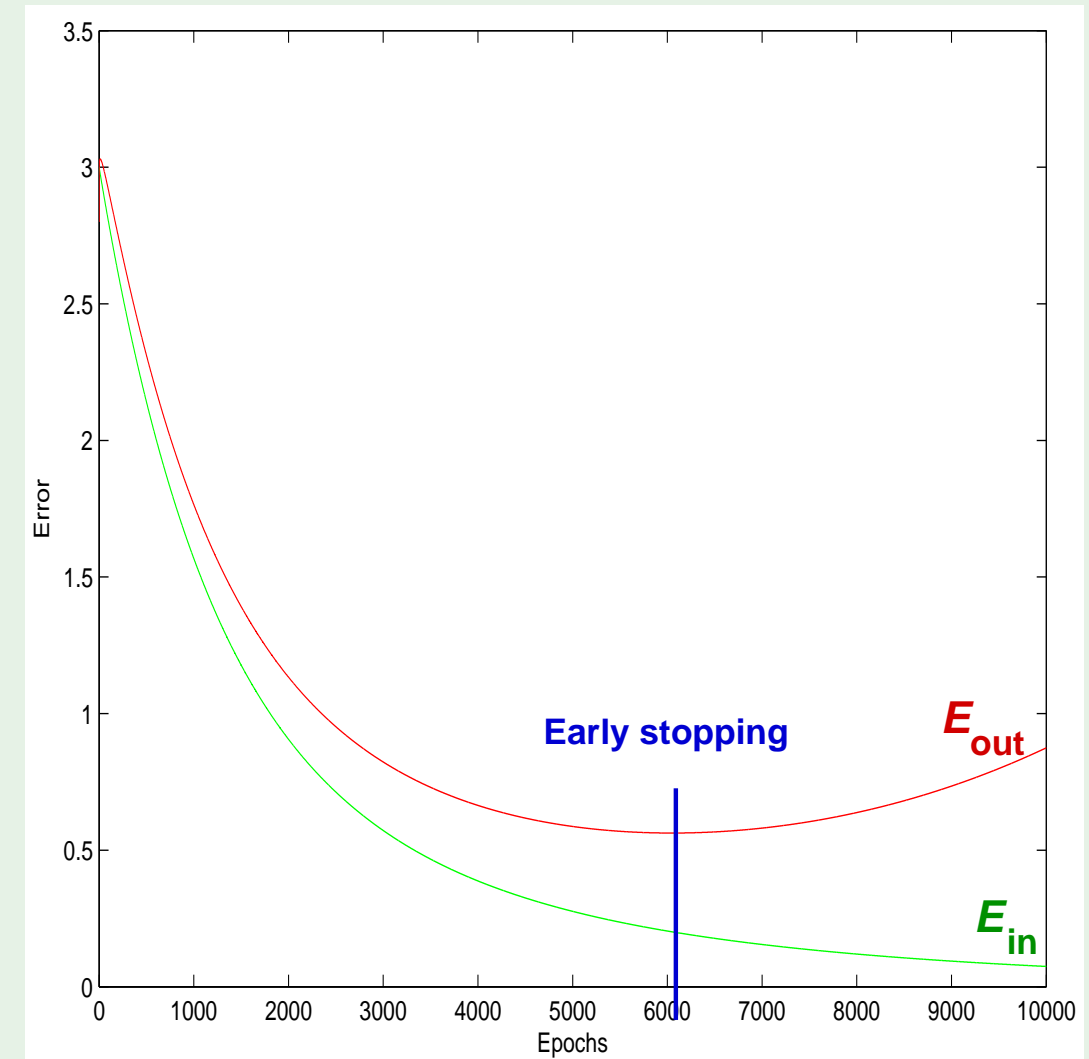**Weight elimination:**

Fewer weights $\implies$ smaller VC dimension



Soft weight elimination:

$$\Omega(\mathbf{w}) = \sum_{i,j,l} \frac{\left(w_{ij}^{(l)}\right)^2}{\beta^2 + \left(w_{ij}^{(l)}\right)^2}$$
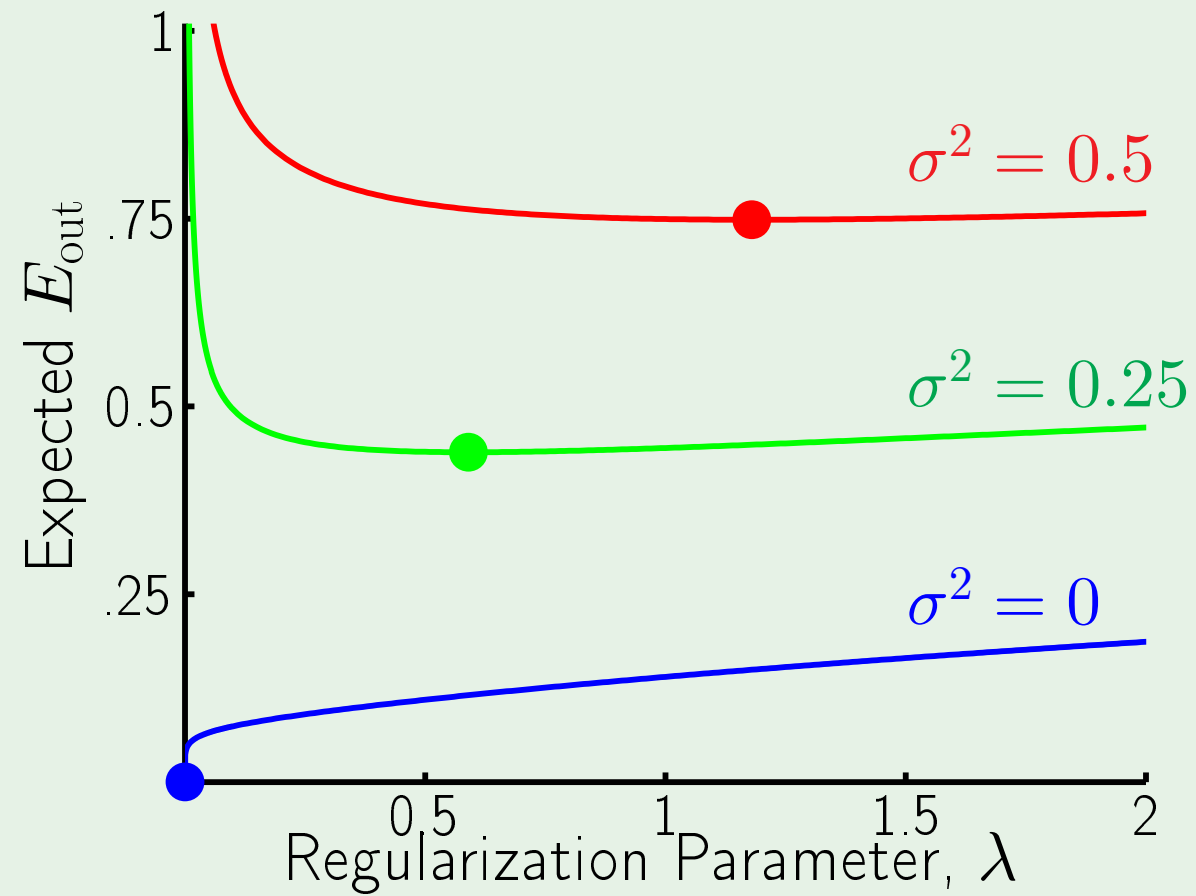
# Early stopping as a regularizer

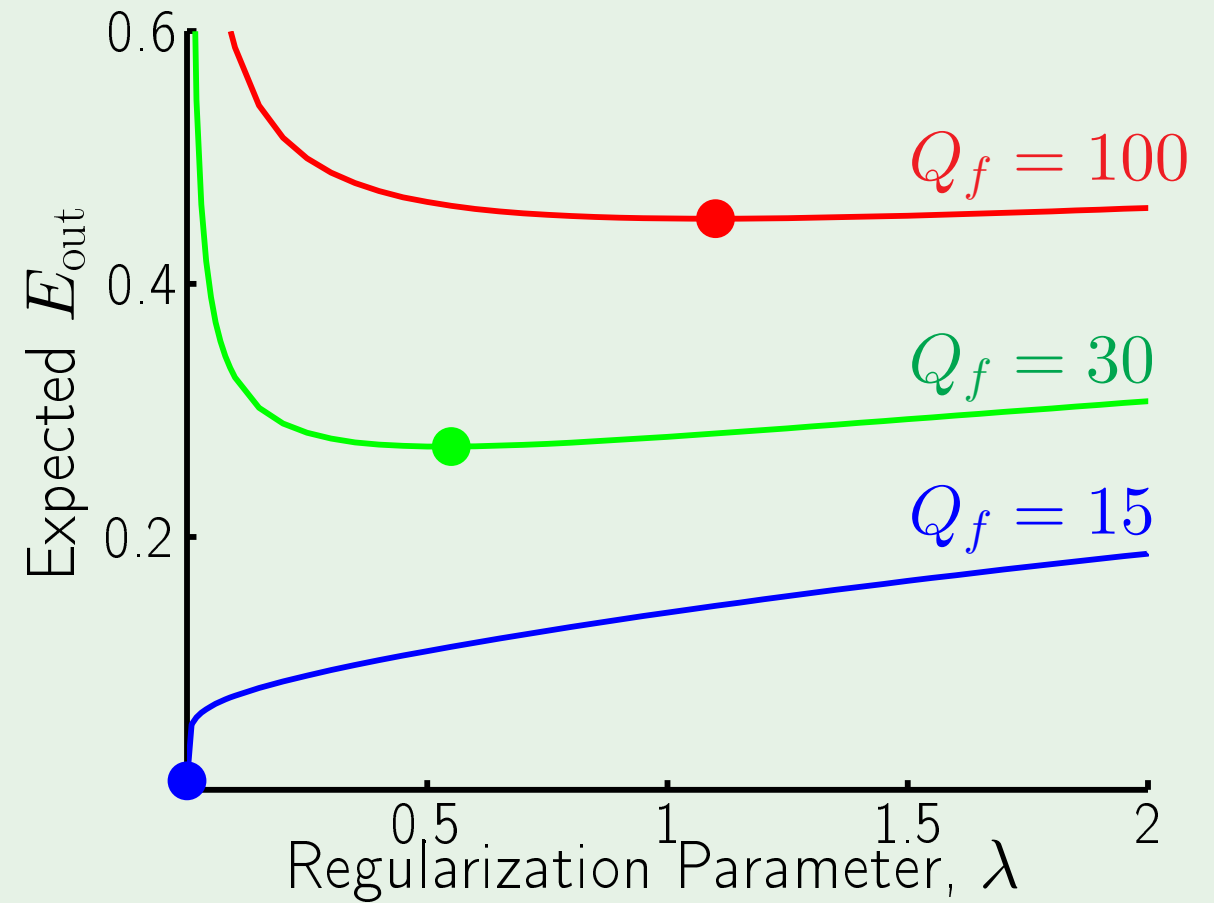Regularization through the optimizer!

When to stop?     **validation**

# The optimal $\lambda$



Stochastic noise

Deterministic noise