# Machine Learning from Data
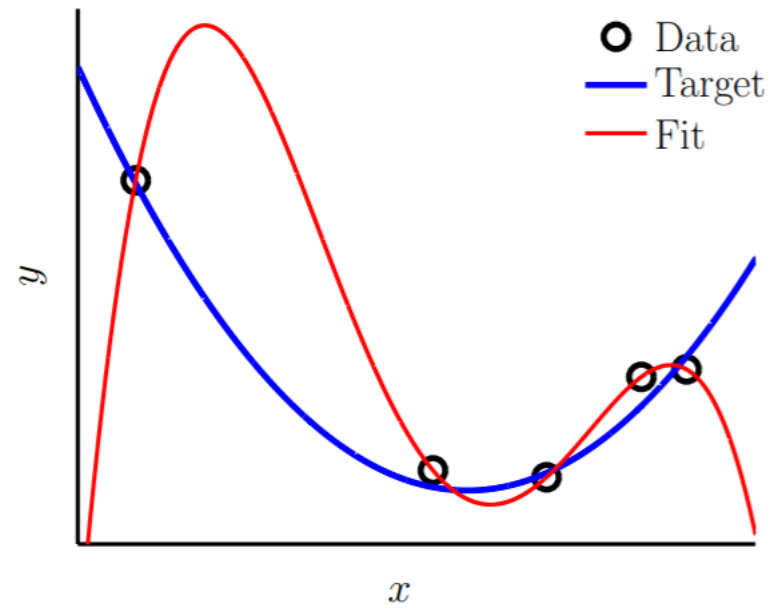
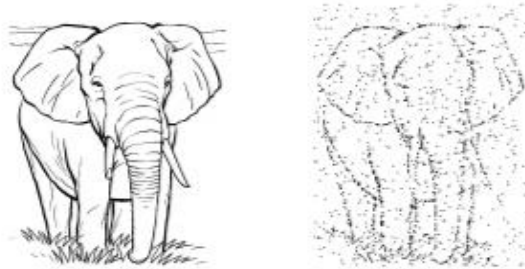Lecture 12: Spring 2021

# Today's Lecture
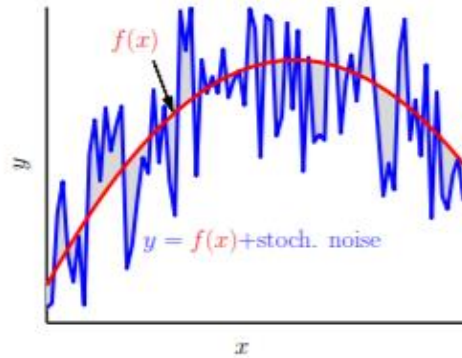
- Regularization
- Constraining the Model
- Augmented Error

# Overfit (Recap)

Fitting the data more than is warranted
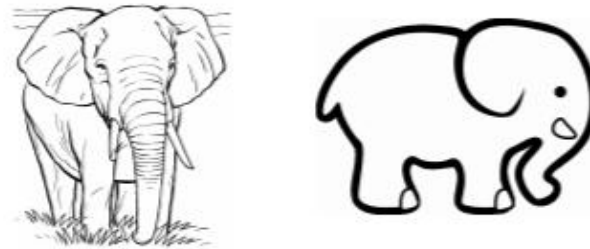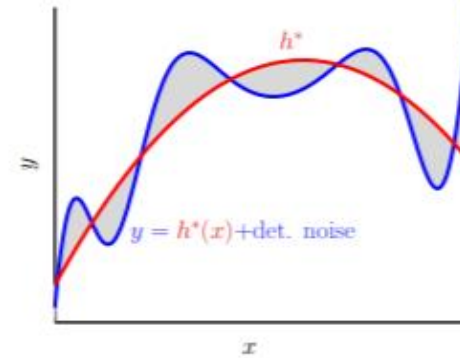
## Stochastic Noise

$f(x)$

$y = f(x) + \text{stoch. noise}$

## Deterministic Noise
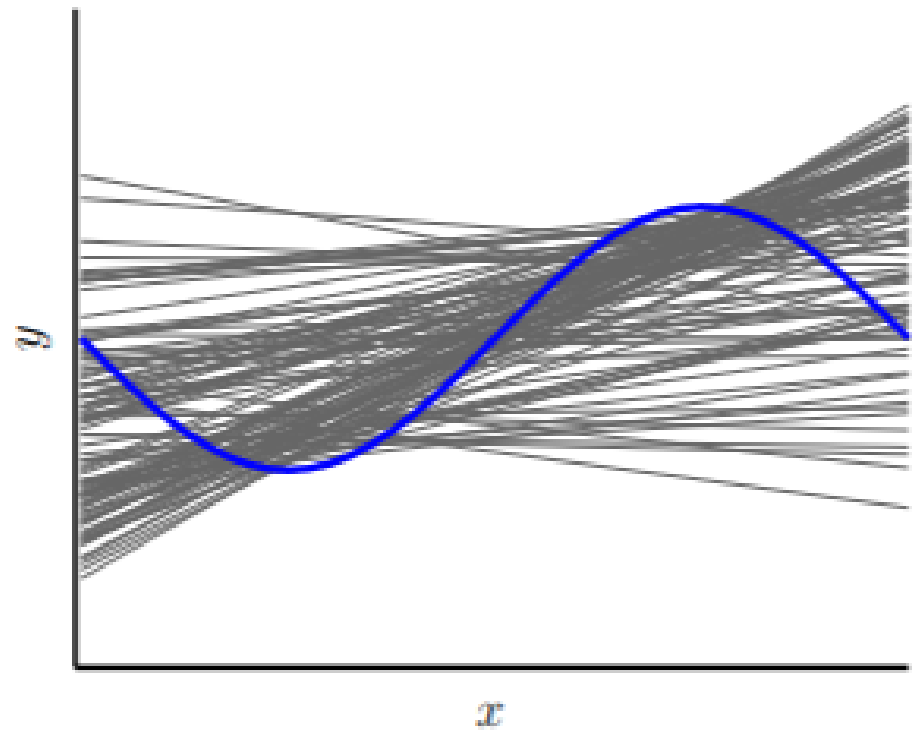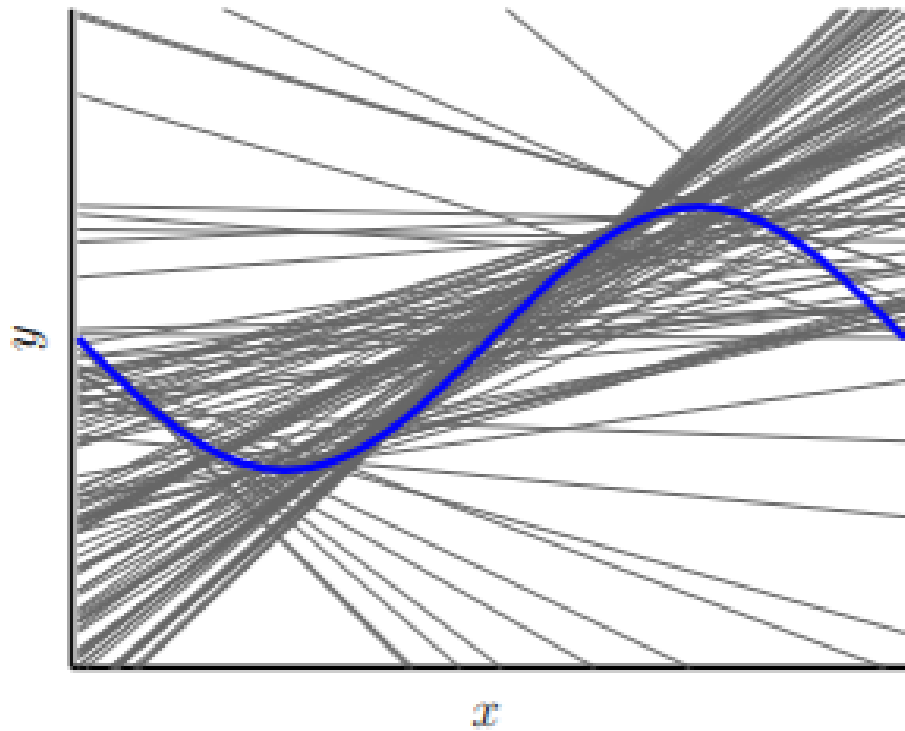
$h^*$

$y = h^*(x) + \text{det. noise}$

## Stochastic and Deterministic Noise Hurt Learning

**Human:** Good at extracting the <u>simple</u> pattern, ignoring the noise and complications.

**Computer:** Pays equal attention to all pixels. Needs help simplifying $\rightarrow$ (features, regularization).
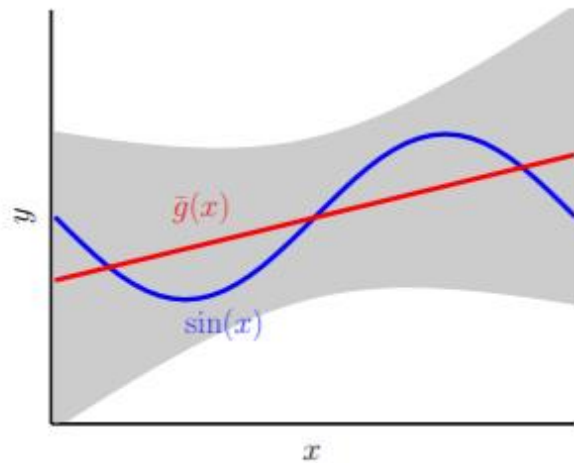
# What is Regularization?

- A cure for our tendency to fit noise, hence improve out-of-sample Error.

- It works by constraining the model so that we cannot fit noise.

- Side effects: If we cannot fit noise maybe we cannot fit the actual signal (f)
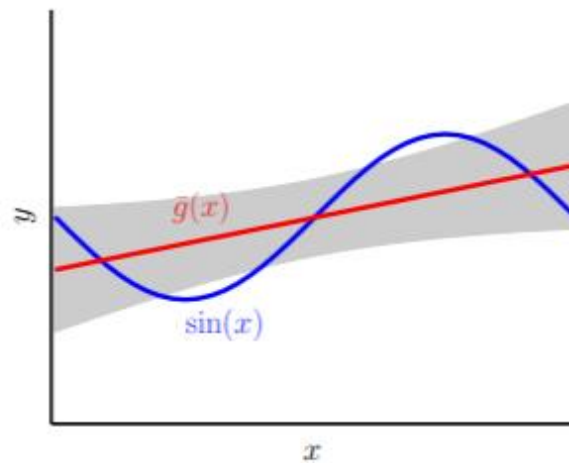
constrain weights to be smaller

Constraining the Model

no regularization

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

regularization

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

← side effect
← treatment

(Constant model had bias=0.5 and var=0.25.)

Bias Variance

# Mathematics of Regularization

-

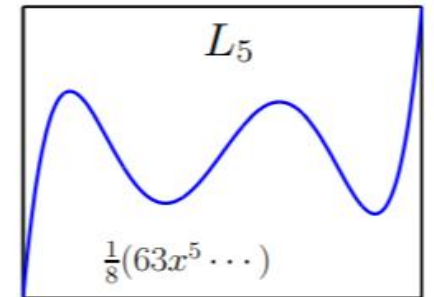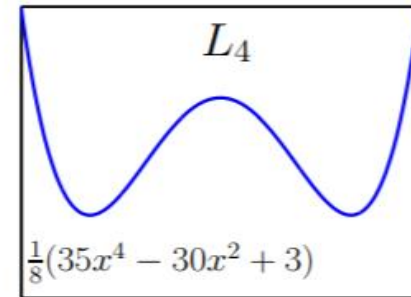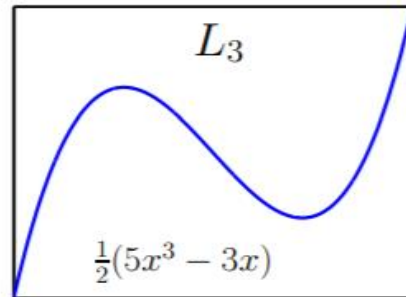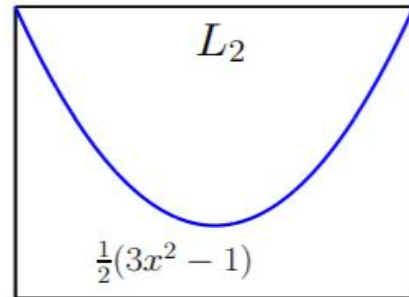$\mathcal{H}_Q$: polynomials of order $Q$.

Standard Polynomial

$$\mathbf{z} = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^Q \end{bmatrix}$$

$$h(x) = \mathbf{w}^{\mathrm{T}}\mathbf{z}(x)$$
$$= w_0 + w_1 x + \cdots + w_Q x^Q$$

Legendre Polynomial

$$\mathbf{z} = \begin{bmatrix} 1 \\ L_1(x) \\ L_2(x) \\ \vdots \\ L_Q(x) \end{bmatrix}$$

we're using linear regression

$$h(x) = \mathbf{w}^{\mathrm{T}}\mathbf{z}(x)$$
$$= w_0 + w_1 L_1(x) + \cdots + w_Q L_Q(x)$$

allows us to treat the weights 'independently'



$L_1$ — $x$

$L_2$ — $\frac{1}{2}(3x^2 - 1)$

$L_3$ — $\frac{1}{2}(5x^3 - 3x)$

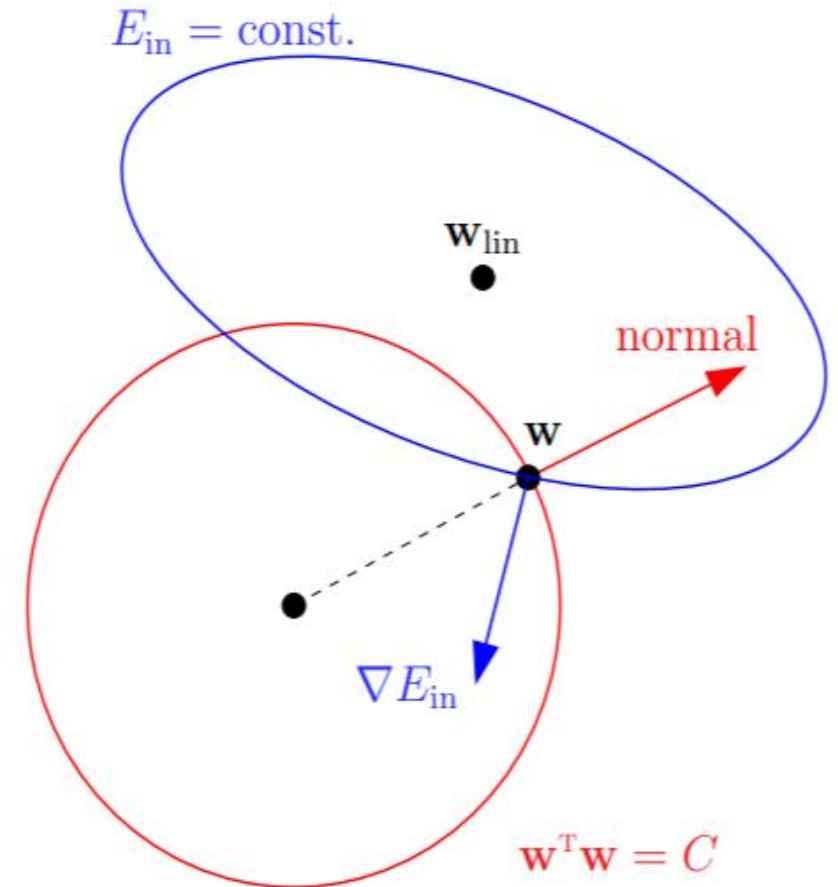$L_4$ — $\frac{1}{8}(35x^4 - 30x^2 + 3)$

$L_5$ — $\frac{1}{8}(63x^5 \cdots)$

$$\text{min}: \quad E_{\text{in}}(\mathbf{w}) = \frac{1}{N}(Z\mathbf{w} - \mathbf{y})^{\text{T}}(Z\mathbf{w} - \mathbf{y})$$

$$\text{subject to:} \quad \mathbf{w}^{\text{T}}\mathbf{w} \leq C$$
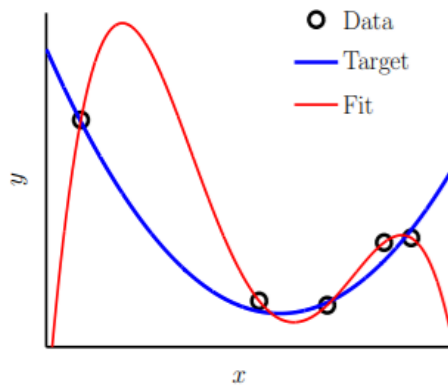
**Observations:**

1. Optimal $\mathbf{w}$ tries to get as 'close' to $\mathbf{w}_{\text{lin}}$ as possible.

   Optimal $\mathbf{w}$ will use full budget and be on the surface $\mathbf{w}^{\text{T}}\mathbf{w} = C$.

2. Surface $\mathbf{w}^{\text{T}}\mathbf{w} = C$, at optimal $\mathbf{w}$, should be perpindicular to $\nabla E_{\text{in}}$.

   Otherwise can move along the surface and decrease $E_{\text{in}}$.

3. Normal to surface $\mathbf{w}^{\text{T}}\mathbf{w} = C$ is the vector $\mathbf{w}$.

4. Surface is $\perp \nabla E_{\text{in}}$; surface is $\perp$ normal.

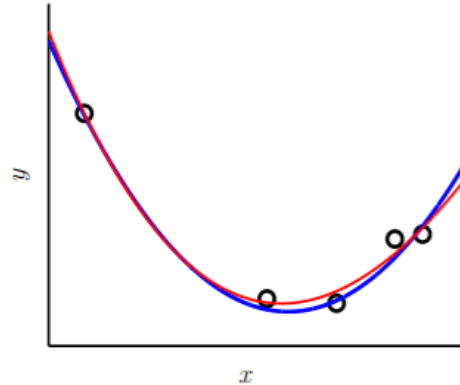   $\nabla E_{\text{in}}$ is parallel to normal (but in opposite direction).

# Regularization In Action

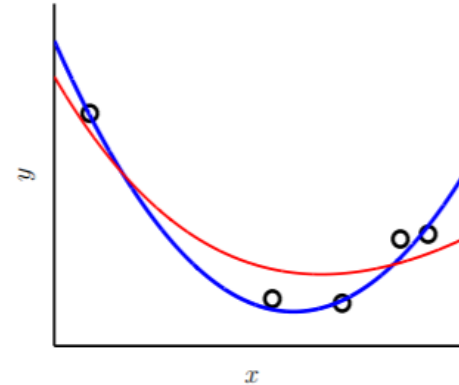Minimizing $E_{\text{in}}(\mathbf{w}) + \dfrac{\lambda}{N}\mathbf{w}^{\text{T}}\mathbf{w}$ with different $\lambda$'s
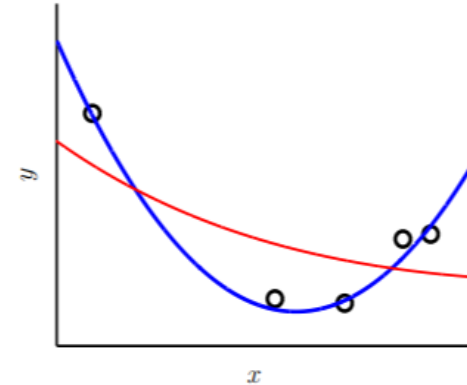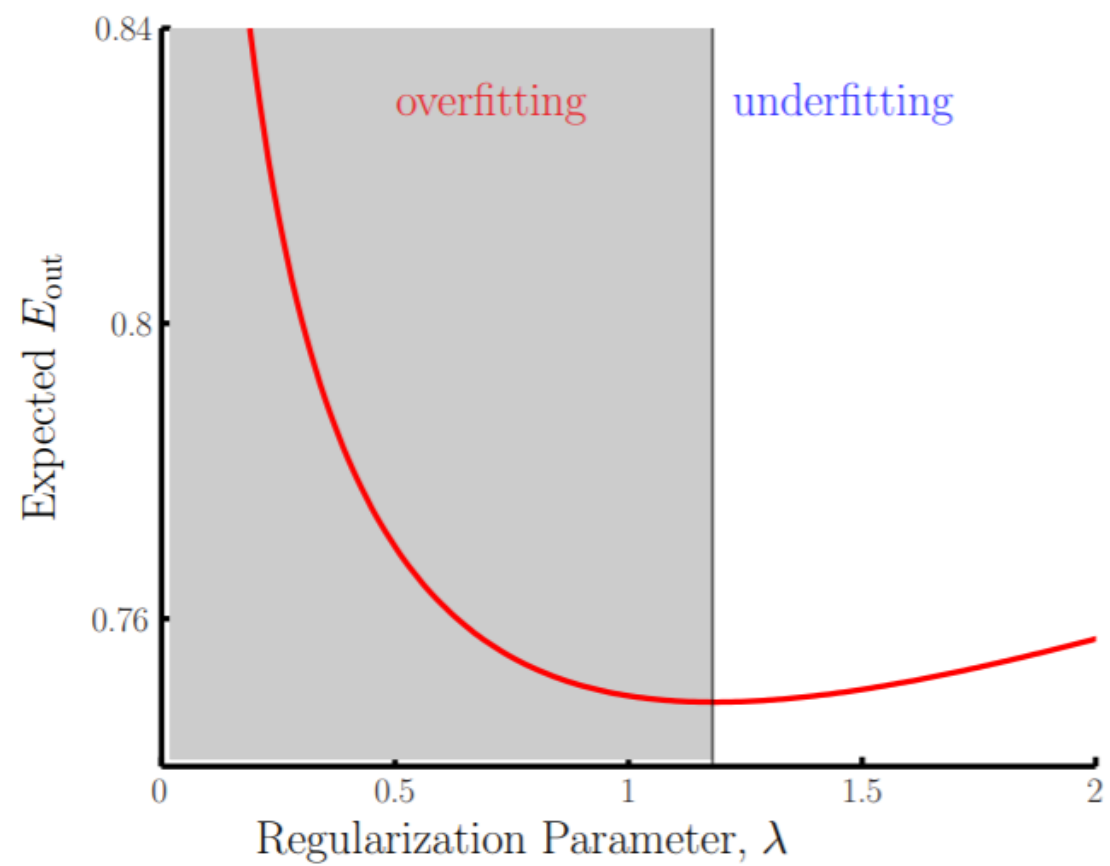
$\underline{\lambda = 0}$      $\underline{\lambda = 0.0001}$      $\underline{\lambda = 0.01}$      $\underline{\lambda = 1}$
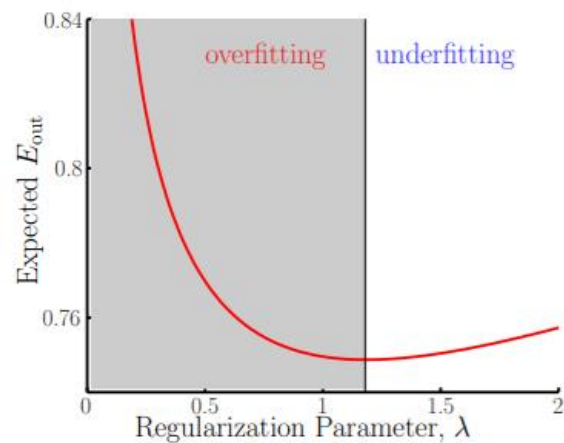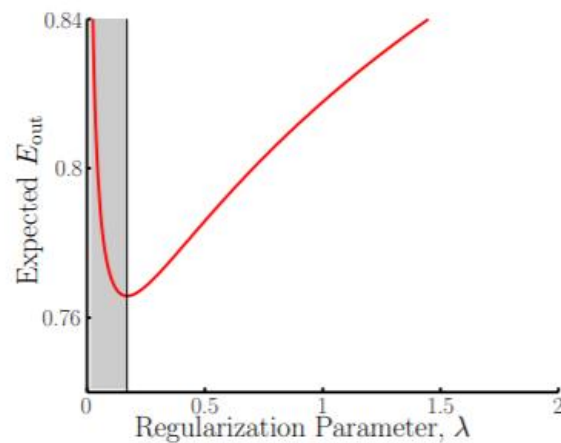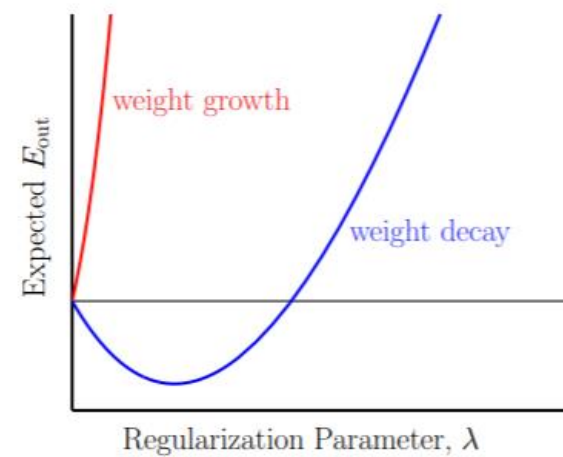
# Uniform Weight Decay



$$\sum_{q=0}^{Q} w_q^2$$

# Low Order Fit



$$\sum_{q=0}^{Q} q w_q^2$$

# Weight Growth!



$$\sum_{q=0}^{Q} \frac{1}{w_q^2}$$

# Choosing a Regularizer – A Practitioner's Guide

The perfect regularizer:

    constrain in the 'direction' of the target function.

    target function is <u>unknown</u> (going around in circles ☺ ).

The guiding principle:

    constrain in the 'direction' of **smoother** (usually simpler) hypotheses

    hurts your ability to fit the 'high frequency' noise

    smoother and simpler $\xrightarrow{\text{usually means}}$ weight decay not weight growth.

Stochastic noise $\longrightarrow$ nothing you can do about that.

Good features $\longrightarrow$ helps to reduce deterministic noise.

Regularization:

Helps to combat what noise remains, especially when $N$ is small.

Typical modus operandi: sacrifice a little **bias** for a huge improvement in **var**.

VC angle: you are using a smaller $\mathcal{H}$ without sacrificing too much $E_{\text{in}}$

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

this was $\mathbf{w}^{\mathsf{T}}\mathbf{w}$

$\updownarrow$

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

this was $O\left(\sqrt{\tfrac{d_{\text{vc}}}{N}\ln N}\right)$

$E_{\text{aug}}$ **can beat** $E_{\text{in}}$ **as a proxy for** $E_{\text{out}}$.

depends on choice of $\lambda$

# Thanks!