




# Machine Learning from Data

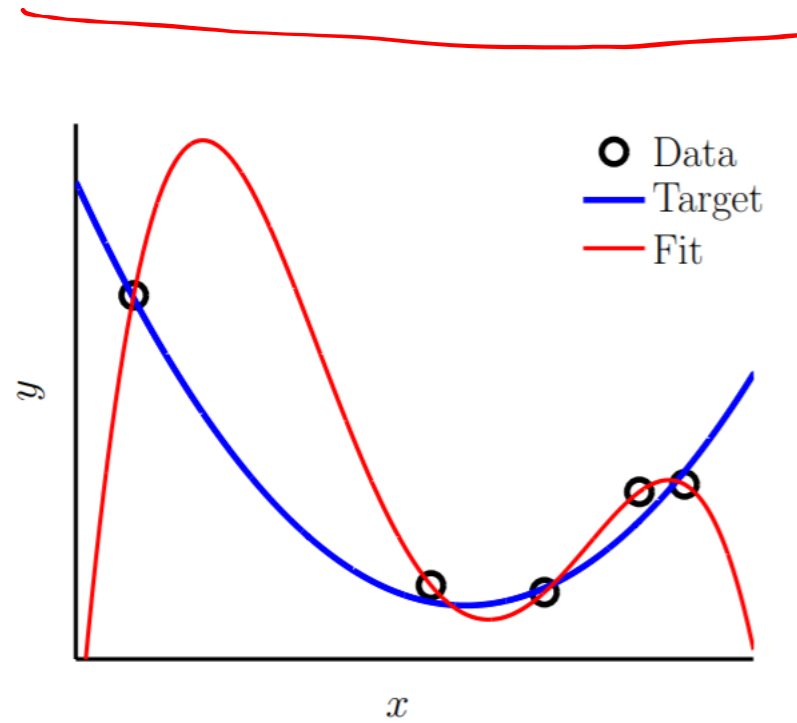
Lecture 12: Spring 2021

# Today's Lecture

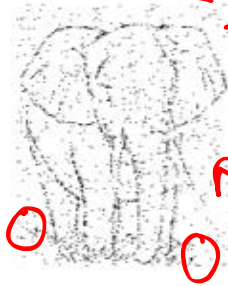
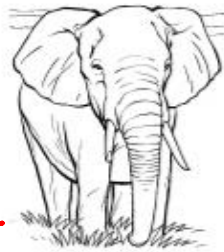
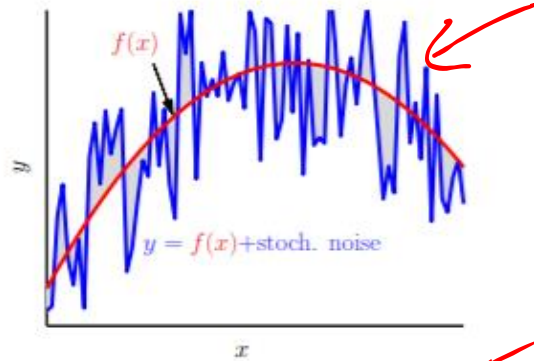
- Regularization 
- Constraining the Model 
- Augmented Error 

# Overfit (Recap)

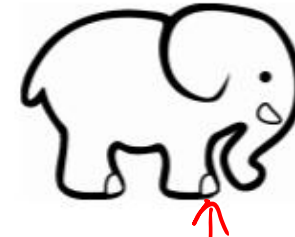
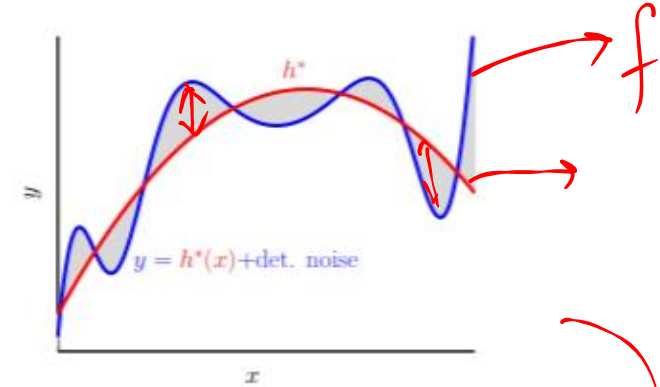
Fitting the data more than is warranted



## Stochastic Noise



## Deterministic Noise



## Stochastic and Deterministic Noise Hurt Learning

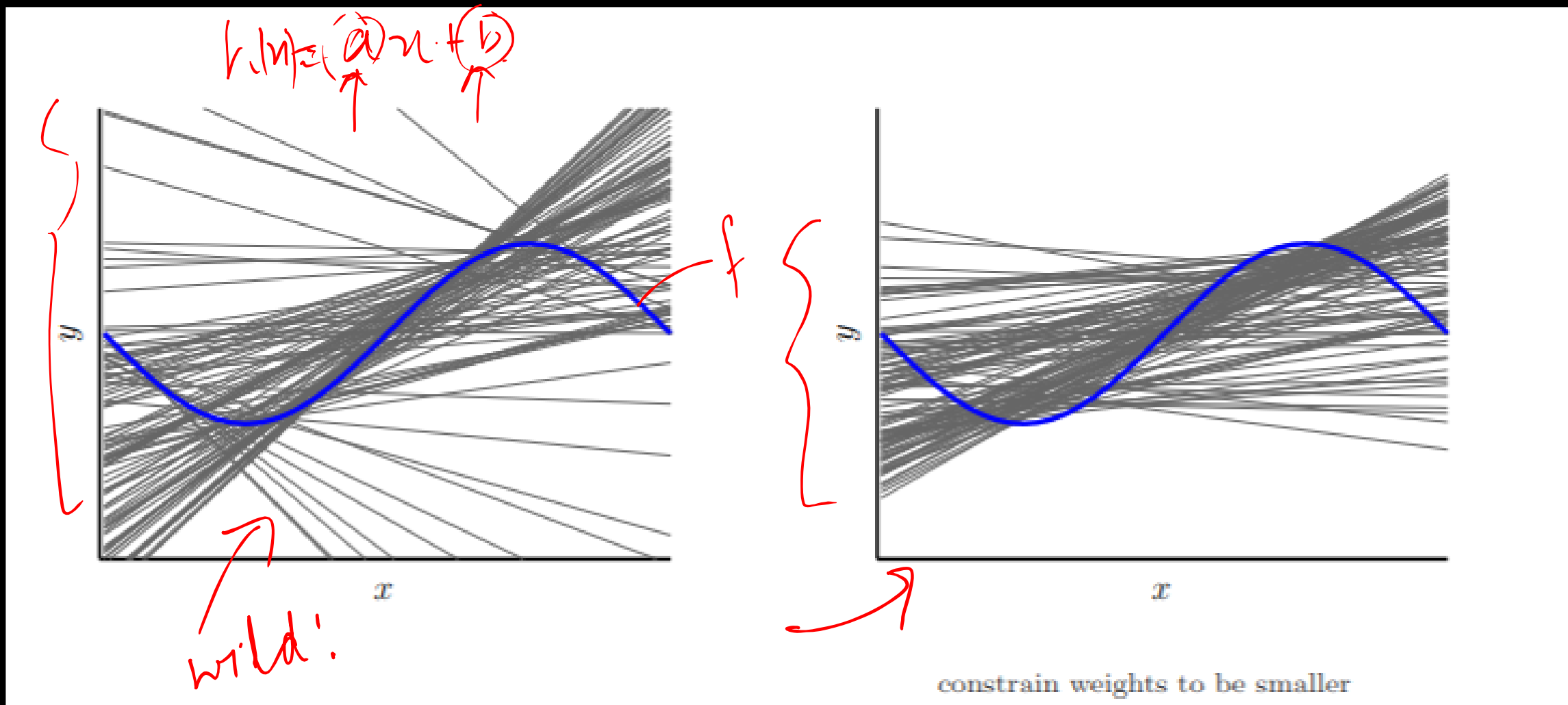
**Human:** Good at extracting the simple pattern, ignoring the noise and complications.

**Computer:** Pays equal attention to all pixels. Needs help simplifying → (features, regularization).

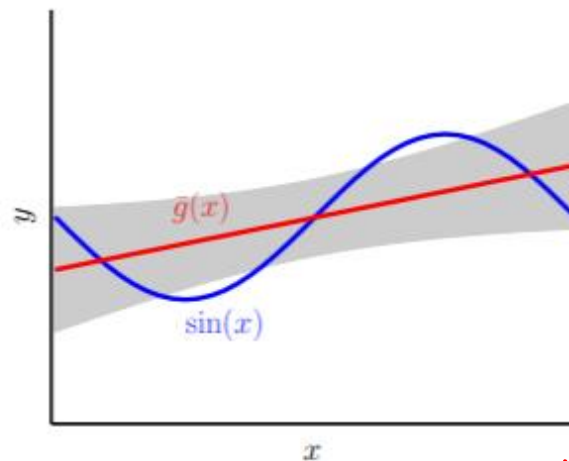
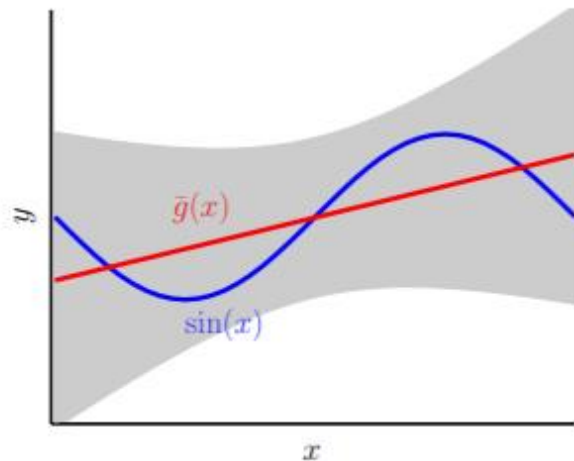
# What is Regularization?

Overfitting  $\rightarrow$  disease

- A cure for our tendency to fit noise, hence improve out-of-sample Error.  $E_{out}$
- It works by constraining the model so that we cannot fit noise. *Medicine*
- Side effects: If we cannot fit noise maybe we cannot fit the actual signal (f)



Constraining the Model



no regularization

bias = 0.21

var = 1.69

regularization

bias = 0.23

var = 0.33

*paying the price*

← side effect

← treatment

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

$$h_b(x) = b$$

$$h(x) = ax + b$$

Bias Variance

# Mathematics of Regularization

- *Constrain the model*

*Linear Models*

*Constraints* → *Hard*  
→ *Soft*

*VC Inequality:*

$$E_{\text{out}}(g) \leq E_{\text{in}}(g) + \underbrace{\Omega(H)}_{\text{higher}}$$

*fitting with a  
simple  $h$*

$\Omega(g)$



$$(x_1, y_1) (x_2, y_2) \dots (x_N, y_N) \xrightarrow{\Phi} (z_1, y_1) (z_2, y_2) \dots (z_N, y_N)$$

Data matrix =  $X$

Target vector =  $y$

Transformed data matrix =  $Z$

find  $\tilde{w} \rightarrow Z$  space

$$\text{Min. } (Z\tilde{w} - y)^T (Z\tilde{w} - y)$$

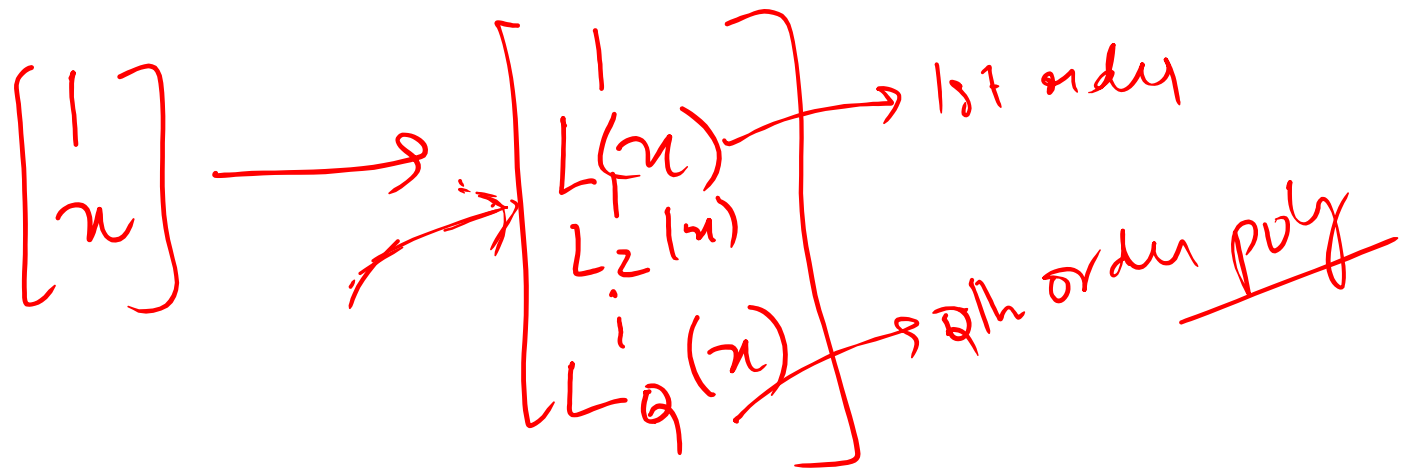
If  $Z \rightarrow$  was invertible

$$\tilde{w} = Z^+ y$$

$$Z\tilde{w} = y$$

$$(Z^+ = (Z^T Z)^{-1} Z^T)$$

Legendre's Polynomials.



$$n^3, n^5$$

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

Standard Polynomial

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

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

Legendre Polynomial

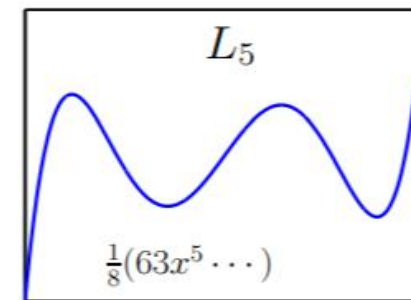
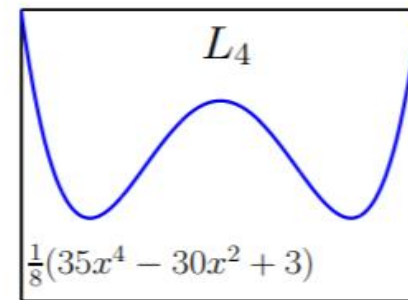
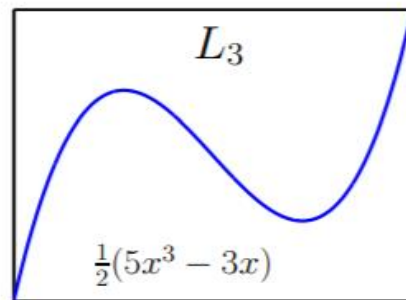
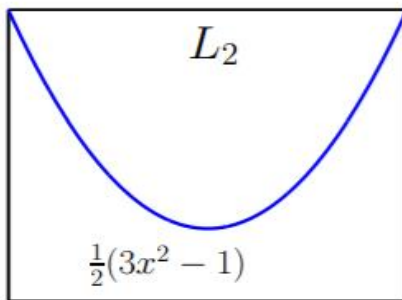
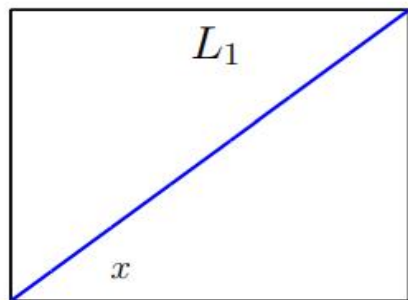
$$\mathbf{z} = \begin{bmatrix} 1 \\ L_1(x) \\ L_2(x) \\ \vdots \\ L_Q(x) \end{bmatrix}$$

we're using linear regression

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

allows us to treat the weights 'independently'

orthogonal



$$\mathcal{H}_2 = \{h(x) \mid h(x) = \omega_0 + \omega_1 l_1(x) + \omega_2 l_2(x), \omega \in \mathbb{R}^3\}$$

$$\mathcal{H}_{10} = \left\{ h(x) \mid h(x) = \omega_0 + \omega_1 l_1(x) + \omega_2 l_2(x) + \dots + \omega_{10} l_{10}(x), \omega \in \mathbb{R}^{11} \right\}$$

$$\mathcal{H}_Q = \sum_{q=0}^Q \omega_q l_q(x)$$

$$\mathcal{H}_2 \subset \mathcal{H}_{10}$$

if  $E_{in}(\mathcal{H}_2) \approx E_{in}(\mathcal{H}_{10})$   
 VC Theory. ↑  
constrain

$$\mathcal{H}_2 = \left\{ h(x) \mid h(x) = \omega_0 + \omega_1 l_1(x) + \omega_2 l_2(x) + \dots + \omega_{10} l_{10}(x), \omega \in \mathbb{R}^{11} \right\}$$

✓

s.t.  $\rightarrow \omega_3 = \omega_4 = \dots = \omega_{10} = 0$

Hard-order constraint.

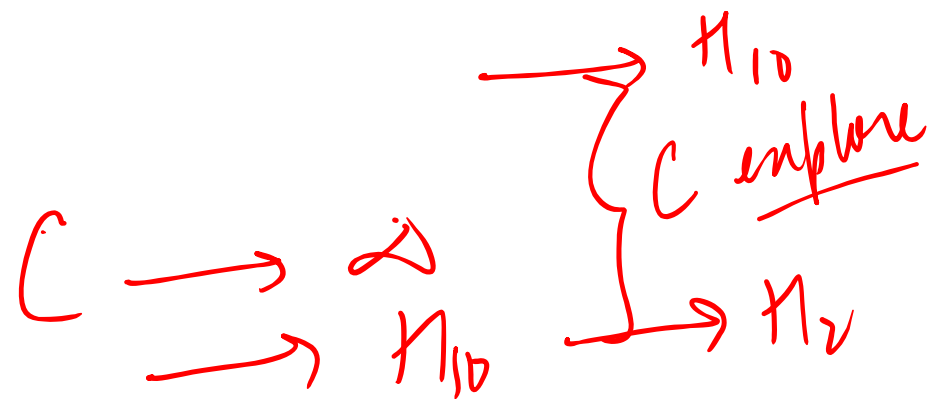
Soft order constraint.

Budget = C

$$\mathcal{H}_C = \left\{ h(x) \mid h(x) = \omega_0 + \omega_1 l_1(x) + \dots + \omega_{10} l_{10}(x), \omega \in \mathbb{R}^{11} \right\}$$

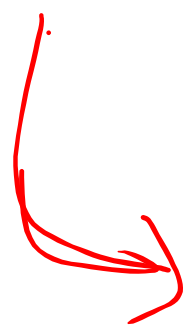
$$\sum_{q=0}^{10} \omega_q^2 \leq C$$

$$\mathcal{H}_C \subset \mathcal{H}_{10}$$



$\tilde{w}$  $\dots N$  $w^T z$  $w^T w \leq C \checkmark$ 

$$E_{in}(w) = \frac{1}{N} (zw - y)^T (zw - y)$$



$$w_{lin} = Z^T y = (Z^T Z)^{-1} Z^T y$$

$$\min_w \underline{E_{in}(w)}$$

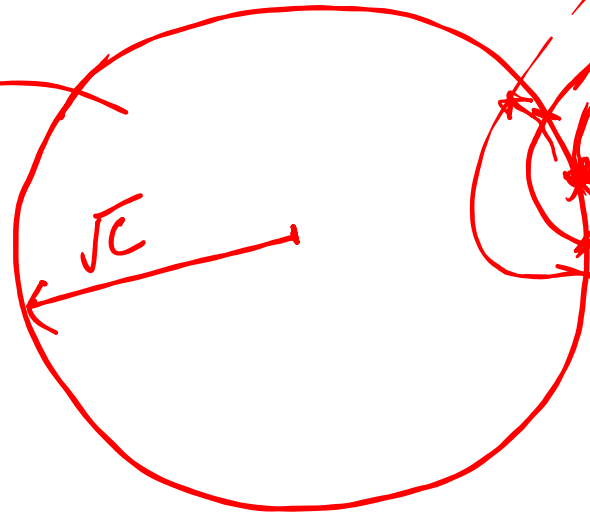
subject to  $w^T w \leq C$

Case 1:  $w_{lin}^T w_{lin} \leq C$

$$w_{reg_\tau} = w_{lin}$$

Case 2:  $w_{lin}^T w_{lin} > C$  ✓

feasible set



$\sqrt{C}$



$E_{in}$  is convex quadratic

$w_{lin}$

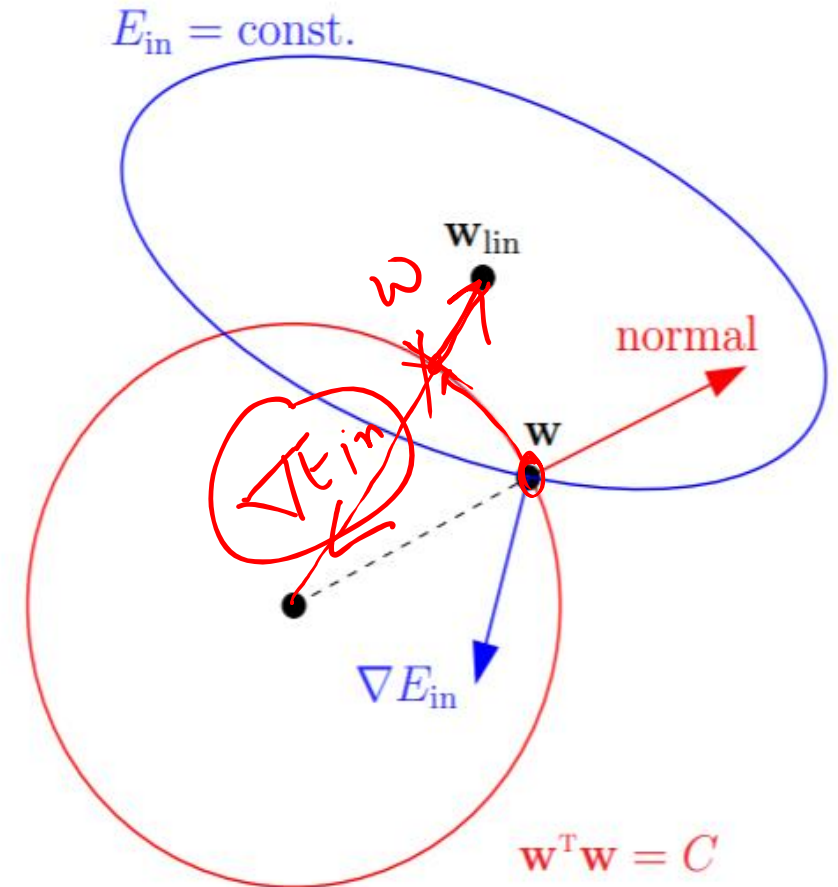
$\Delta E_{in}$

$$\begin{aligned} \min : \quad & E_{\text{in}}(\mathbf{w}) = \frac{1}{N}(\mathbf{Z}\mathbf{w} - \mathbf{y})^T(\mathbf{Z}\mathbf{w} - \mathbf{y}) \\ \text{subject to:} \quad & \mathbf{w}^T \mathbf{w} \leq C \end{aligned}$$

$$\nabla E_{\text{in}} \propto -\mathbf{w}$$

### Observations:

- 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  $\mathbf{w}^T \mathbf{w} = C$ .
- Surface  $\mathbf{w}^T \mathbf{w} = C$ , at optimal  $\mathbf{w}$ , should be perpendicular to  $\nabla E_{\text{in}}$ .  
 Otherwise can move along the surface and decrease  $E_{\text{in}}$ .
- Normal** to surface  $\mathbf{w}^T \mathbf{w} = C$  is the vector  $\mathbf{w}$ .
- Surface is  $\perp \nabla E_{\text{in}}$ ; surface is  $\perp$  **normal**.  
 $\nabla E_{\text{in}}$  is parallel to **normal** (but in opposite direction).





Necessary condition for optimality

$$\nabla E_{in} \propto -w$$

Lagrange's multiplier  $> 0$

$$\text{or } \nabla E_{in} = -2\lambda_c w$$

mathematical convenience

$$\nabla E_{in} + 2\lambda_c w = 0$$

gradient of  $\lambda_c w^T w$

$$\nabla (E_{in} + \lambda_c w^T w) = 0$$

condition.

Constrained optimization



Unconstrained optimization

$$\min E_{in}(w) + \lambda_c w^T w$$

$\lambda_c \uparrow \rightarrow c$  must go down

Fitting under a constraint (Budget)

Unconstrained  $\Downarrow$  penalized fitting (penalty parameter  $\lambda$ )

$$\min E_{\text{avg}}(w) = E_{\text{in}}(w) + \underbrace{\frac{\lambda}{N} w^T w}_{\text{penalty term}}$$

$$E_{\text{avg}}(w) = E_{\text{in}}(w) + \frac{\lambda}{N} \Omega(w)$$

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

complexity of hypothesis.

$w^T w$  — Regularizer (we pick)

$\lambda$  — Regularization

parameter knob control

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

Weight Decay

$E_{\text{avg.}}$

$\underbrace{\lambda}_{\text{weight decay}} (w^T w) \rightarrow$  weight decay.

$w_{\text{reg.}}$

$$\min E_{\text{in}}(w) + \lambda/N w^T w$$

$$\min \frac{1}{N} (Zw - y)^T (Zw - y) + \frac{\lambda}{N} w^T w$$

$$\min w^T Z^T Z w - 2w^T Z^T y + y^T y + \lambda \underline{w^T w}$$

$$\min w^T (Z^T Z + \lambda I) w - 2w^T Z^T y + y^T y$$

gradient  $\rightarrow 0$

$$2 \left( \underbrace{Z^T Z}_{\substack{\text{true} \\ \text{semi definite}}} + \underbrace{\lambda I}_{\text{true definite}} \right) w - 2 Z^T y = 0$$

Invertible

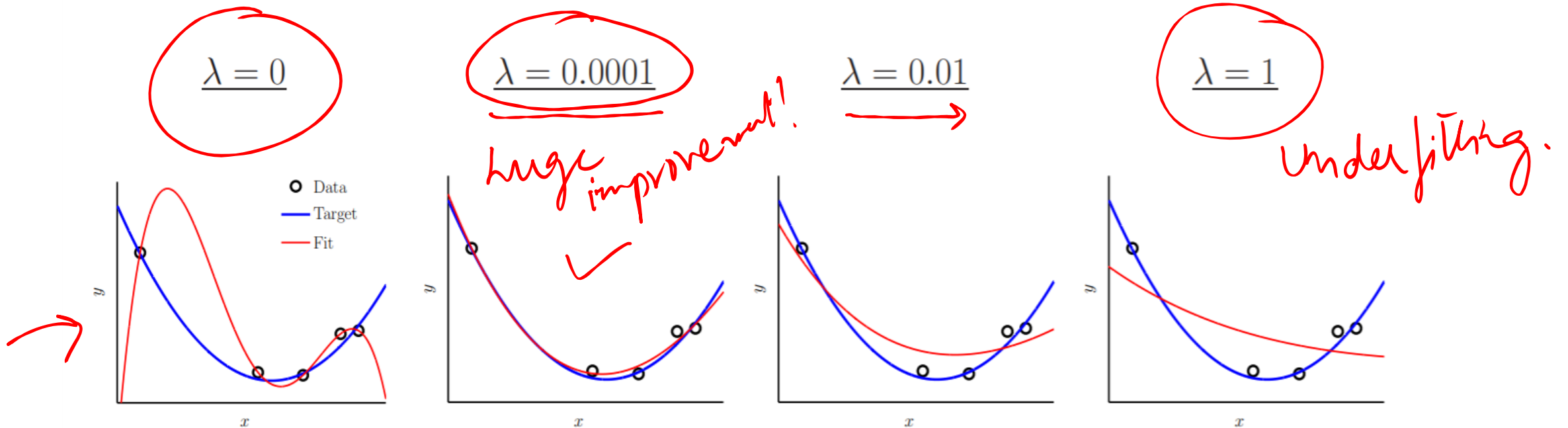
Solve for  $w$

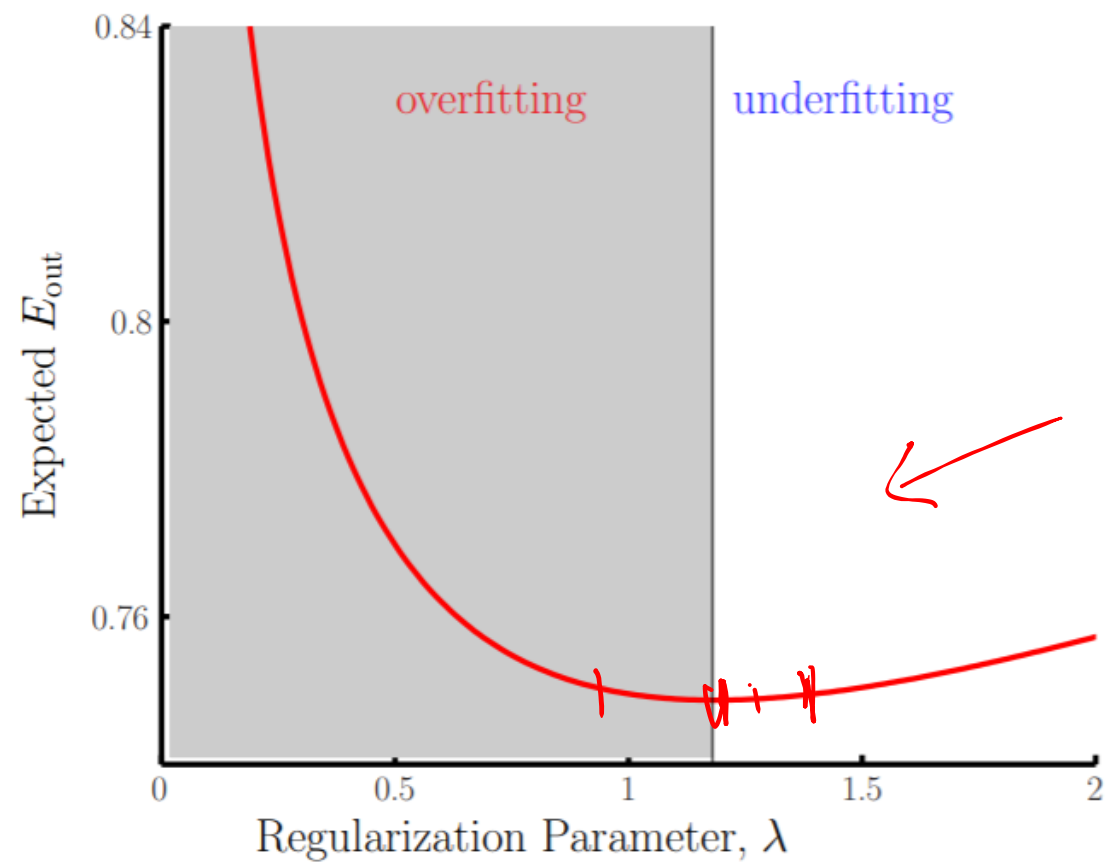
$$w_{\text{reg}} = \underbrace{(Z^T Z + \lambda I)^{-1} Z^T y}_{\text{Regularized linear regression feature transform.}}$$

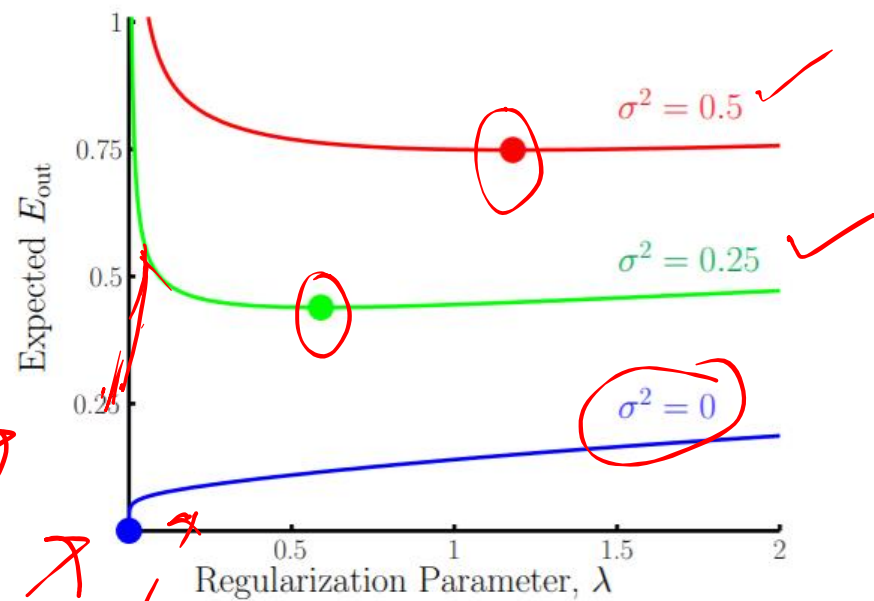
$\lambda \rightarrow$  penalizing large weights.

# Regularization In Action

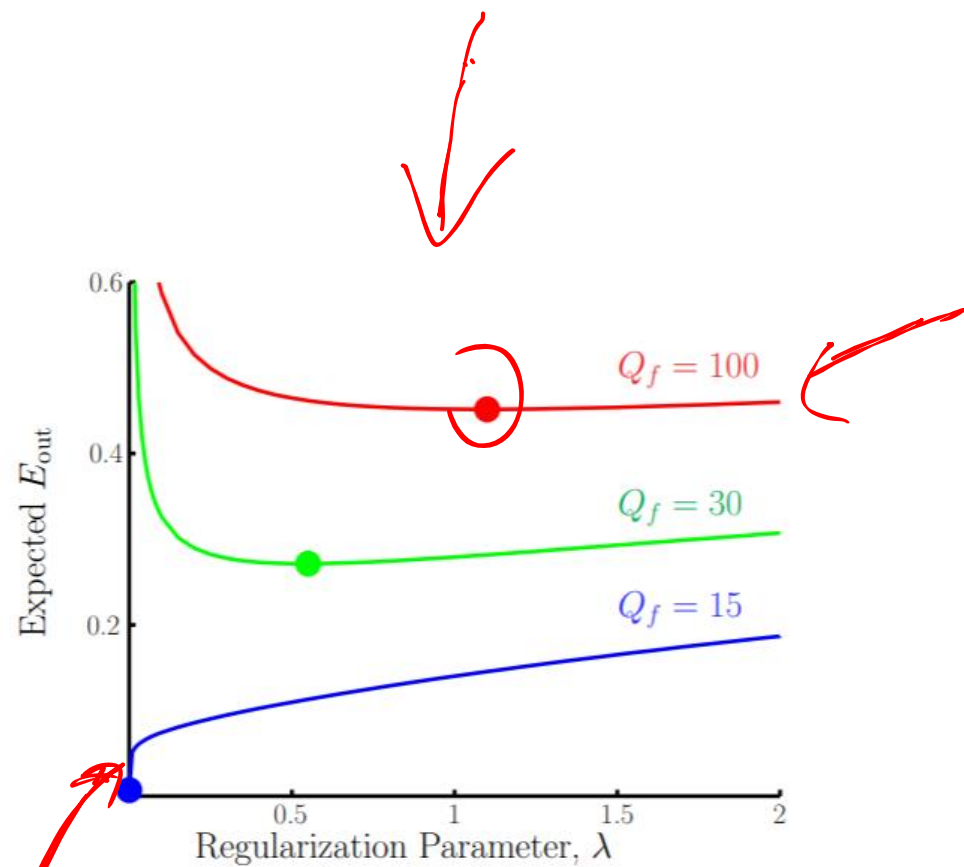
Minimizing  $E_{\text{in}}(\mathbf{w}) + \frac{\lambda}{N} \mathbf{w}^T \mathbf{w}$  with different  $\lambda$ 's





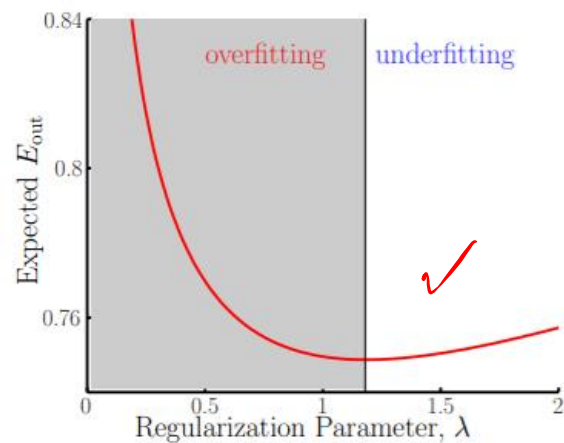


Stochastic

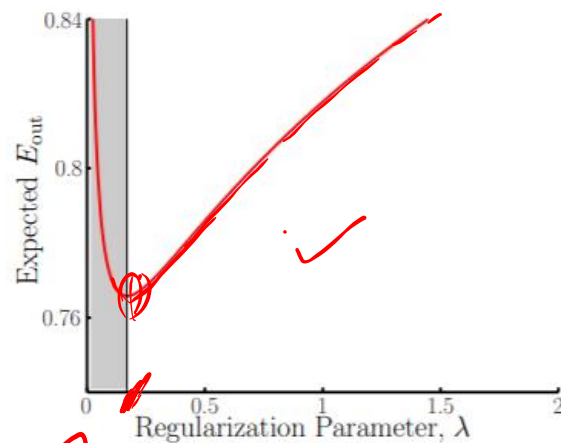


Deterministic

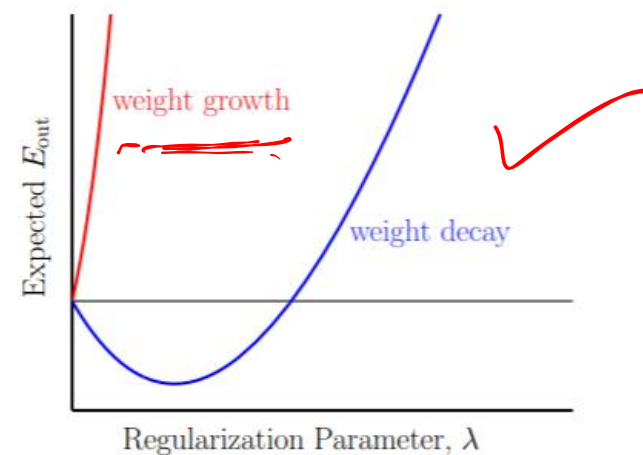
## Uniform Weight Decay



## Low Order Fit



## Weight Growth!



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

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

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

$\Omega_{unit}$



# Choosing a Regularizer – A Practitioner's Guide

The perfect regularizer:

- constrain in the 'direction' of the target function.
- target function is unknown (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  $E_{\text{in}}$

$$\begin{array}{c}
 \underline{E_{\text{aug}}(h)} = \underline{E_{\text{in}}(h)} + \underline{\frac{\lambda}{N}\Omega(h)} \quad \leftarrow \text{this was } \underbrace{\mathbf{w}^T \mathbf{w}} \\
 \Downarrow \\
 \underline{E_{\text{out}}(h)} \leq E_{\text{in}}(h) + \underbrace{\Omega(\mathcal{H})} \quad \leftarrow \text{this was } O\left(\sqrt{\frac{d_{\text{vc}}}{N} \ln N}\right)
 \end{array}$$

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

$\nwarrow$  depends on choice of  $\lambda$

Thanks!