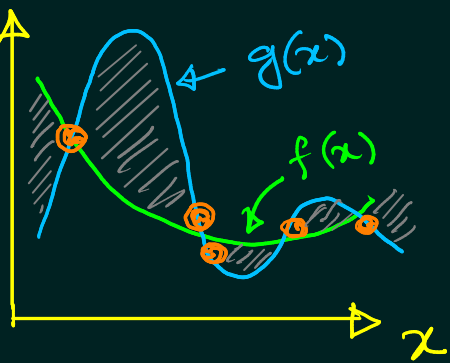


# Overfitting: $y$

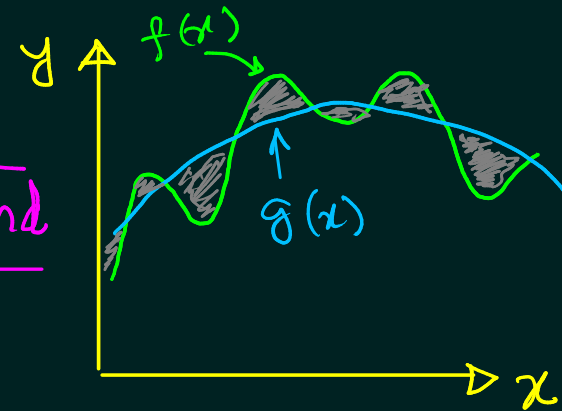
① Stochastic Noise

② Deterministic Noise



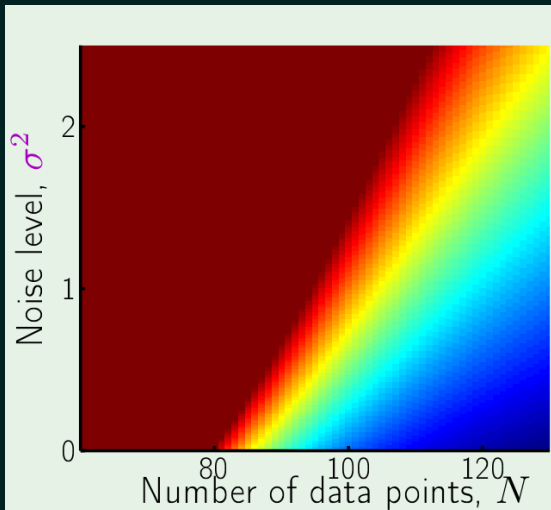
Data with Noise fitting  
①

and

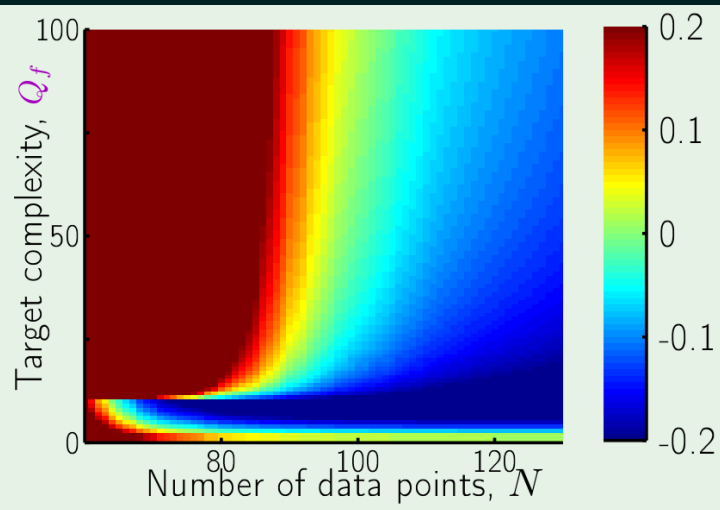


Higher order Data fitting  
②

## SUMMARY



Stochastic noise



Deterministic noise

- ↑ Training data → ↓ Overfitting
- ↑ Noise level → ↑ Overfitting
- ↑ Target Complexity → ↑ Overfitting

→ Relation with VC-Dimension:

↳ Generalization Bound ⇒

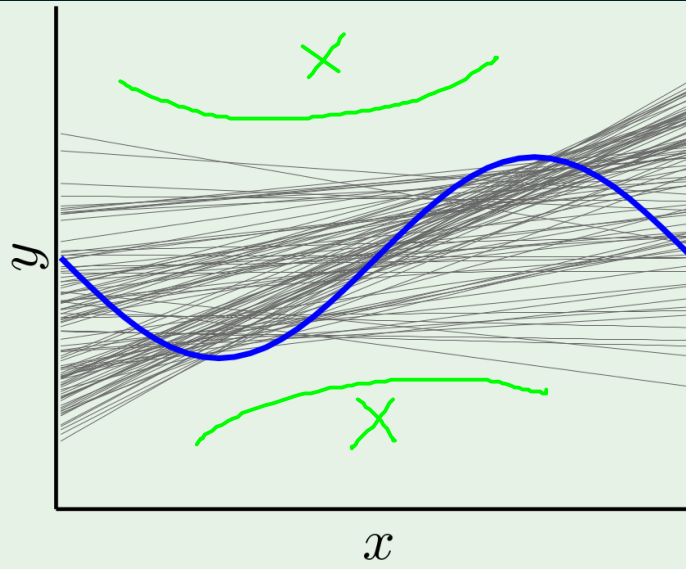
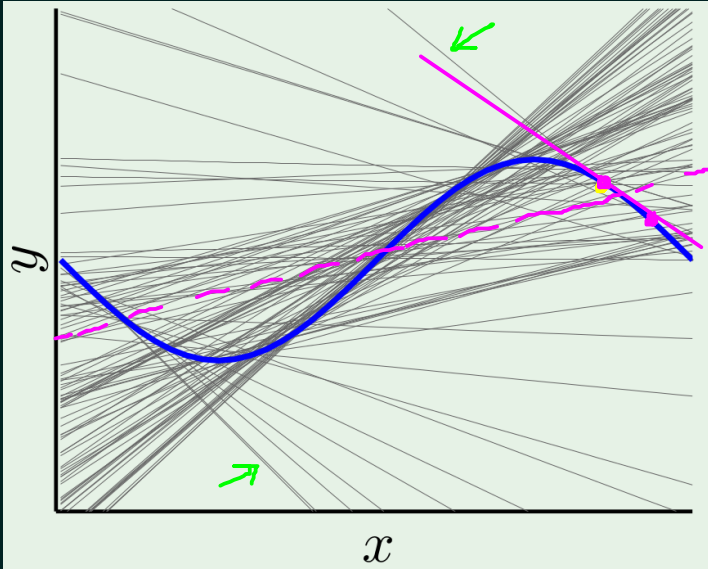
Allows it w but Never Predicts this ??

→ Overfitting Remedy:

① REGULARIZATION + ② VALIDATION

THIS WEEK'S COVERAGE

# Regularization to Rescue

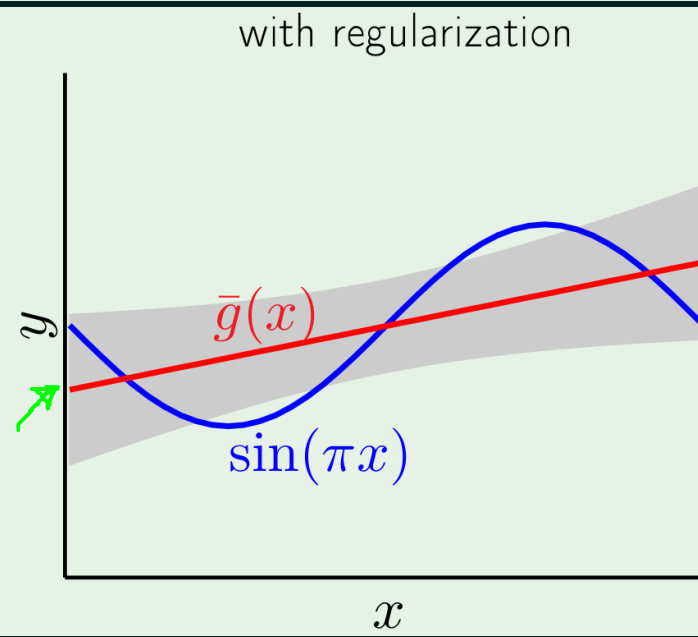
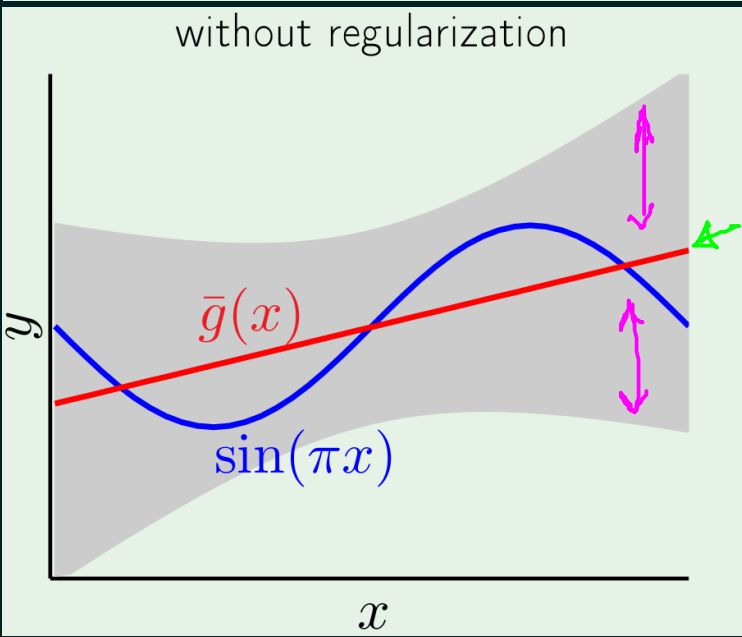


Not Allow  
all slopes

$\checkmark$   $f(x) = ax + b$   
 $\checkmark$

$f(x) = c$

Constraining  
solution



bias = 0.21    var = 1.69

bias = 0.23    var = 0.33

$$\mathcal{H}_Q = \left\{ \sum_{q=0}^Q \underbrace{w_q}_{\checkmark} L_q(x) \right\}$$

$$Z = \begin{bmatrix} 1 \\ L_1(x) \\ \vdots \\ L_Q(x) \end{bmatrix}$$

$$(x_1, y_1) \dots (x_N, y_N)$$

$$\Downarrow$$

$$(z_1, y_1) \dots (z_N, y_N)$$

$$\text{Min. } E_{in}(w) = \frac{1}{N} \sum_{n=1}^N (w^T z_n - y_n)^2$$

$$= \frac{1}{N} (Zw - Y)^T (Zw - Y)$$

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

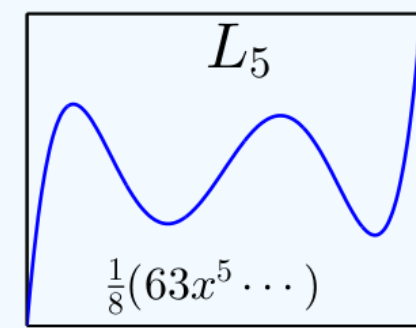
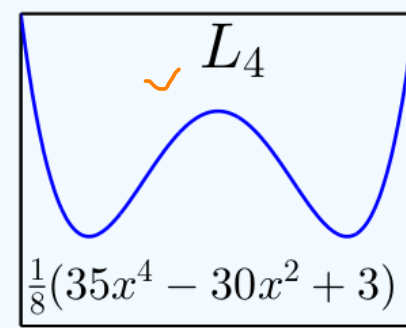
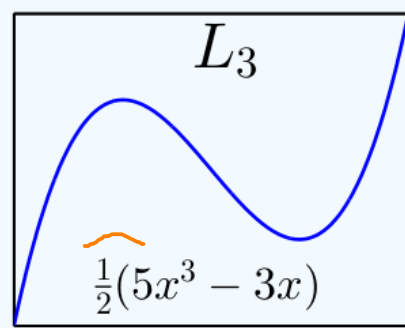
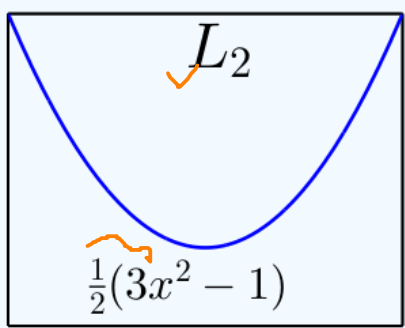
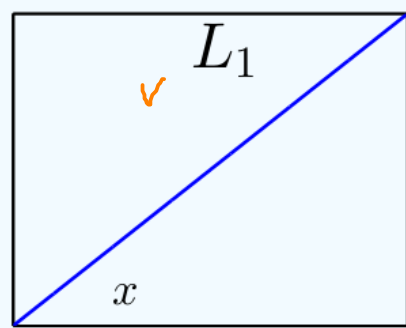
Linear Regression

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

$\mathcal{H}_q$  with  $w_q = 0$  for  $q > 2$   
(Hard)

Unconstrained Fitting of N points

Legendre Polynomials



$$\mathcal{H}_{\frac{1}{2}} = \left\{ \sum_{q=0}^Q w_q L_q(x) \right\}$$

Soft Constraint:  $\sum_{q=0}^Q w_q^2 \leq C$

Minimize  $E_{in}(w) = \frac{1}{N} (Zw - Y)^T (Zw - Y)$   
 subject to  $W^T W \leq C$

Dual

$$\nabla E_{in}(w_{reg}) \propto -w_{reg}$$

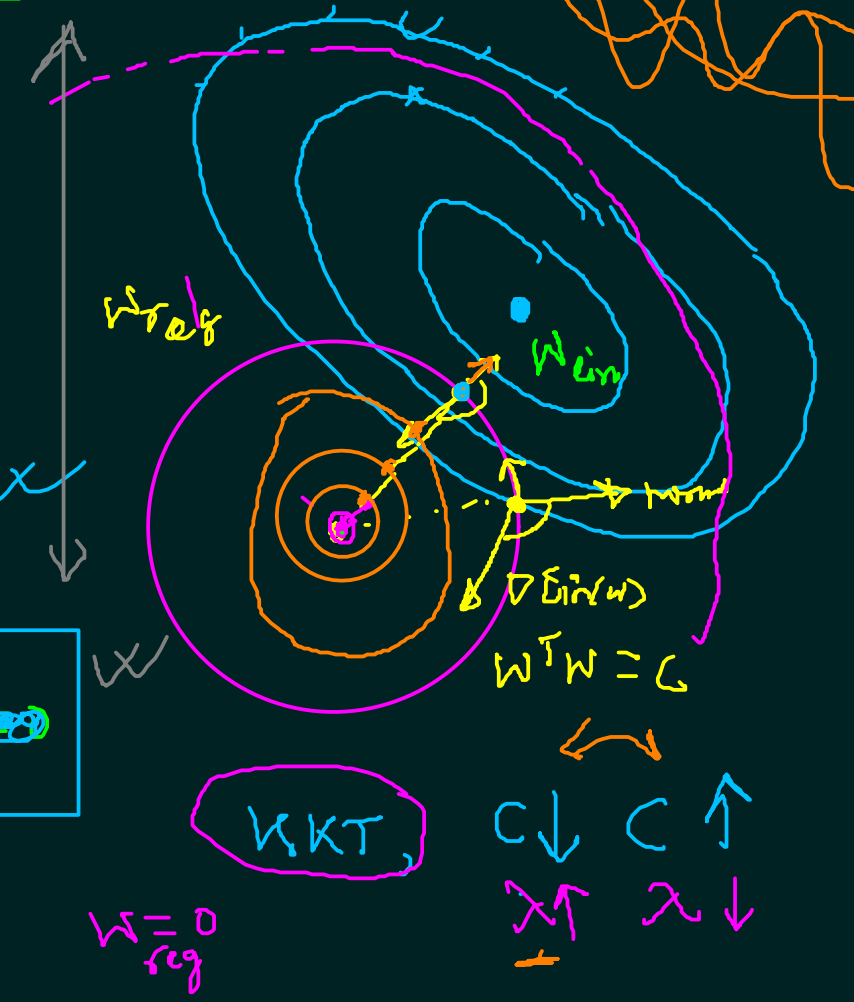
$$= -2 \frac{\lambda}{N} w_{reg}$$

$$\Rightarrow \nabla E_{in}(w_{reg}) + 2 \frac{\lambda}{N} w_{reg} = 0$$

Min.  $E_{in}(w) + \frac{\lambda}{N} W^T W$

st:  $\lambda \geq 0$

$\lambda > 0 \Rightarrow \lambda = 0$



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

$$= \frac{1}{N} \left[ \underbrace{(Zw - y)^T (Zw - y)}_{\text{Q.P. Tech.}} + \lambda w^T w \right]$$

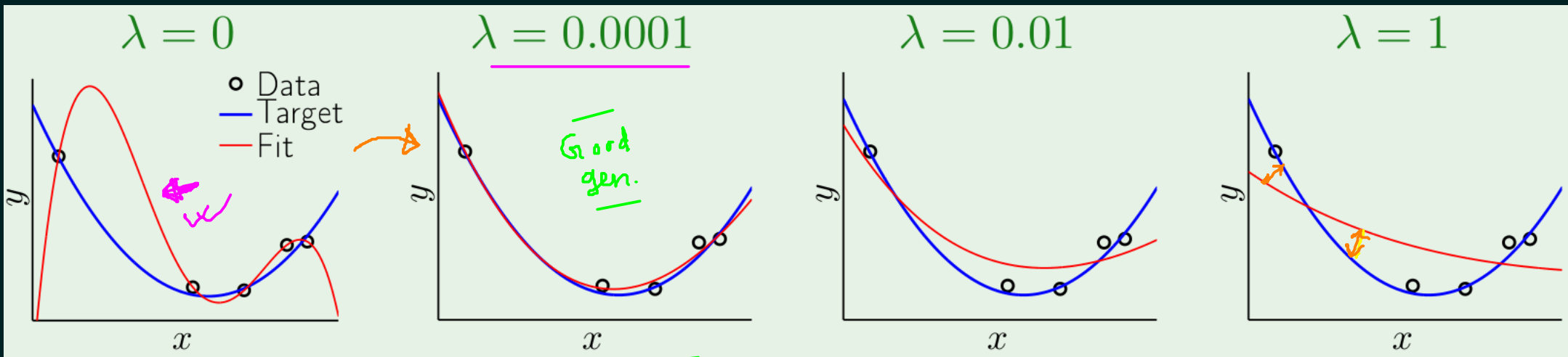
$$\nabla E_{\text{in}}(w) = 0 \Rightarrow Z^T (Zw - y) + \lambda w = 0$$

(Reg)  $\Rightarrow w_{\text{reg}} = \frac{(Z^T Z + \lambda I)^{-1} Z^T y}{}$

Opposed to;  $w_{\text{lin}} = \frac{(Z^T Z)^{-1} Z^T y}{}$  (w/o Reg)

$\lambda \uparrow \text{inf} \Rightarrow w_{\text{reg}} \approx 0$  /  $\lambda \downarrow \approx 0 \Rightarrow w_{\text{reg}} \approx w_{\text{lin}}$

$\rightarrow$  More  $\lambda \Rightarrow$  flatter curves  
(Smoother)

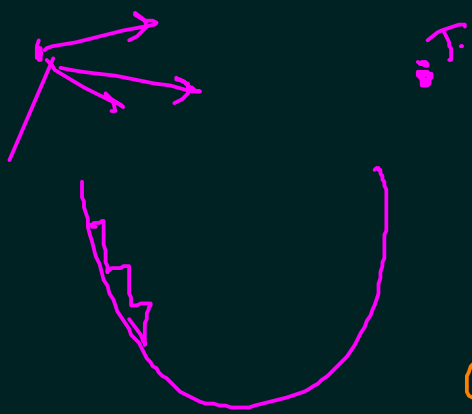


overfitting  $\rightarrow$  Reg 1  $\rightarrow$  Reg' 1  $\rightarrow$  ...  $\rightarrow$  underfitting

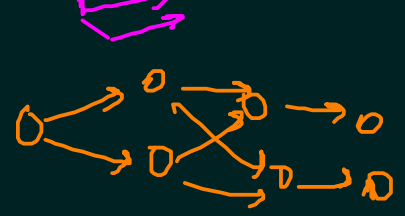
Weight decay:

Gradient descent

$$w(t+1) = w(t) - \eta \left[ \nabla E_{in}(w(t)) + \underbrace{2 \frac{\lambda}{N} w(t)} \right]$$



$$= w(t) \underbrace{\left[ 1 - 2\eta \frac{\lambda}{N} \right]}_{\text{decay}} - \eta \nabla E_{in}(w(t))$$



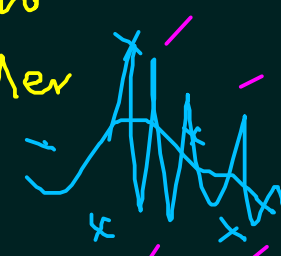
$$\left[ w^T w = \sum_{l=1}^L \sum_{i=0}^{d^{l-1}} \sum_{j=1}^{d^l} (w_{ij}^l)^2 \right]$$

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

$\gamma_q = 2^q \rightarrow$  low order poly  
 $\gamma_q = 2^{-q} \rightarrow$  high order poly

①  $\sum w_q^2 \leq C \rightarrow$  const weight smaller

②  $\sum w_q^2 > C$  ✗



$\rightarrow$  Stochastic Noise "high-freq"

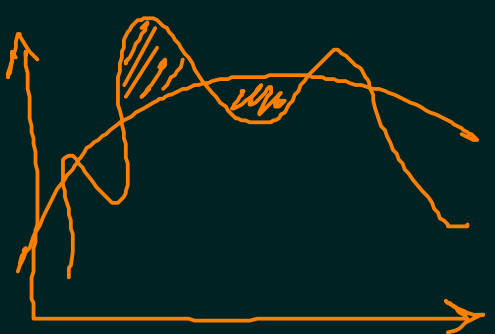
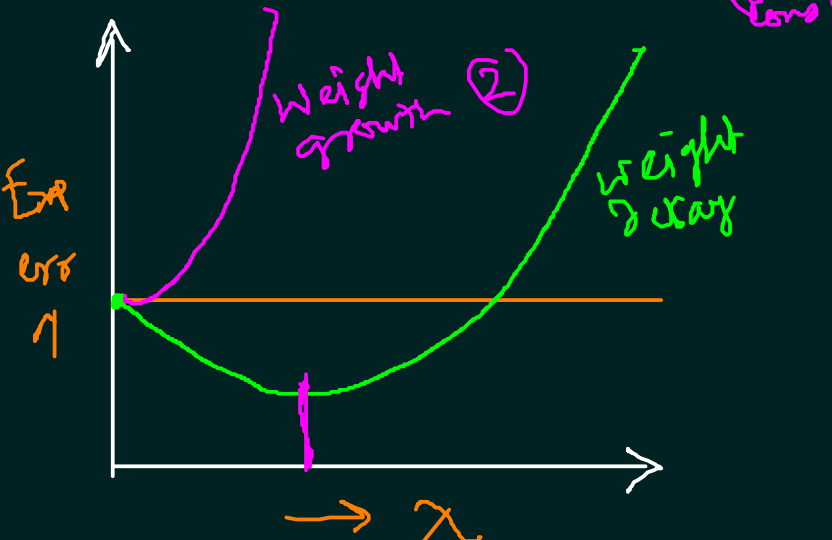
$\rightarrow$  Deterministic Noise "non-smooth"

Practical rule

$\Rightarrow$  "Smother by pathosis"

$$\mathcal{H}_2 > \mathcal{H}_0$$

[OCCAM RAZOR]



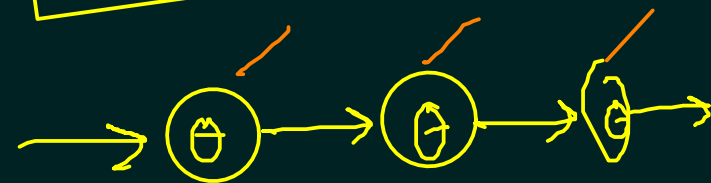
$$E_{\text{avg}}(h) = E_{\text{in}}(h) + \frac{\lambda}{N} \Omega(h)$$

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

$E_{\text{avg}}(h)$  is closer est. / better than  $E_{\text{in}}$  as proxy for  $E_{\text{out}}$



$$w_i x_i$$



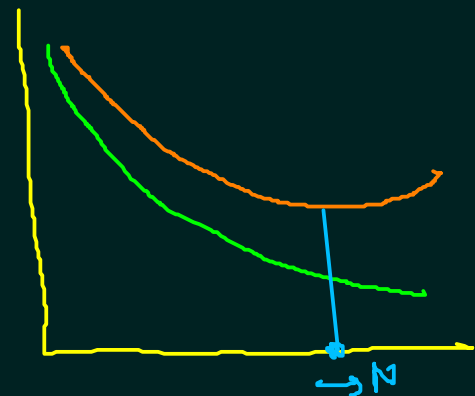
$$\Omega(w) = \sum \frac{(w_{ij}^{(l)})^2}{\beta^2 + (w_{ij}^{(l)})^2}$$

Soft weight elimination

Reg  $\rightarrow$   $\lambda$  (validation)  
 $\rightarrow$   $Z$  (Legendre poly)

validation

$$\mathcal{D}_{\text{train}} = \mathcal{D}_t \cup \mathcal{D}_v$$



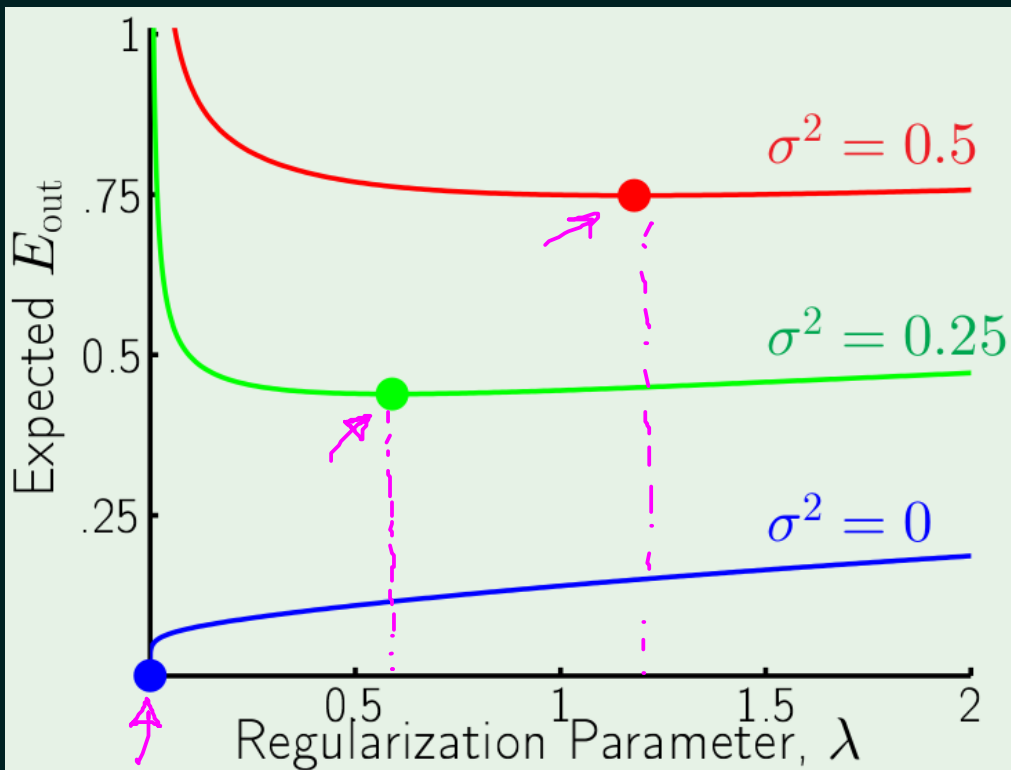


Optimal  $\lambda$

↙ Deterministic }  
↘ Stochastic }

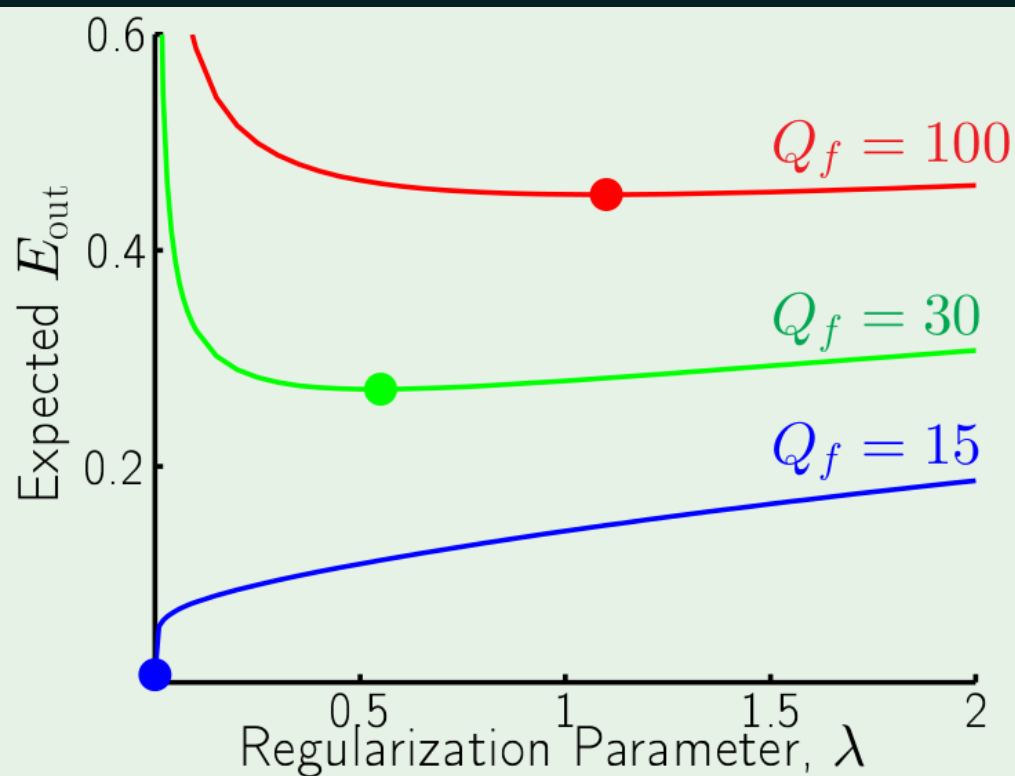
Validation

$\lambda$



Stochastic noise

( $\sigma^2$ )



Deterministic noise

( $Q_f$ )