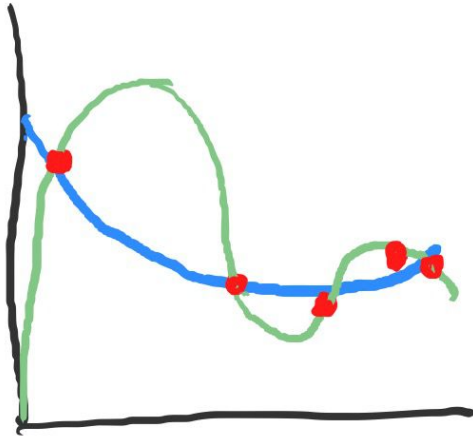
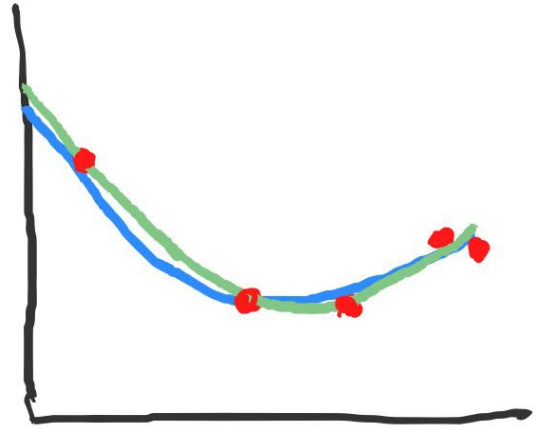


Regularization

Overcoming overfit



(free fit situation)
(No constraint)
(4th order poly)



(Restricted fit)
(4th order poly,
but with a
regularization term)

Regularization : It constraints the learning algorithm to improve out-sample error when noise is present or we have a lack of training samples.

Regularization from VC analysis :

$$E_{\text{out}}(h) \leq E_{\text{in}}(h) + \underbrace{\Omega(\mathcal{H})}_{\text{model complexity penalty}} \quad \forall h \in \mathcal{H}$$

(depends on \mathcal{H}, N, δ)

Idea : Fit using a 'simple' h from \mathcal{H}
so effectively minimizing

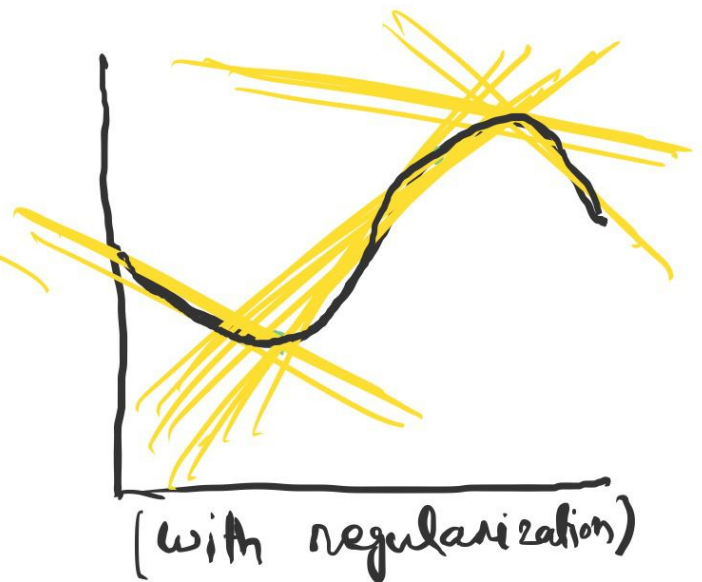
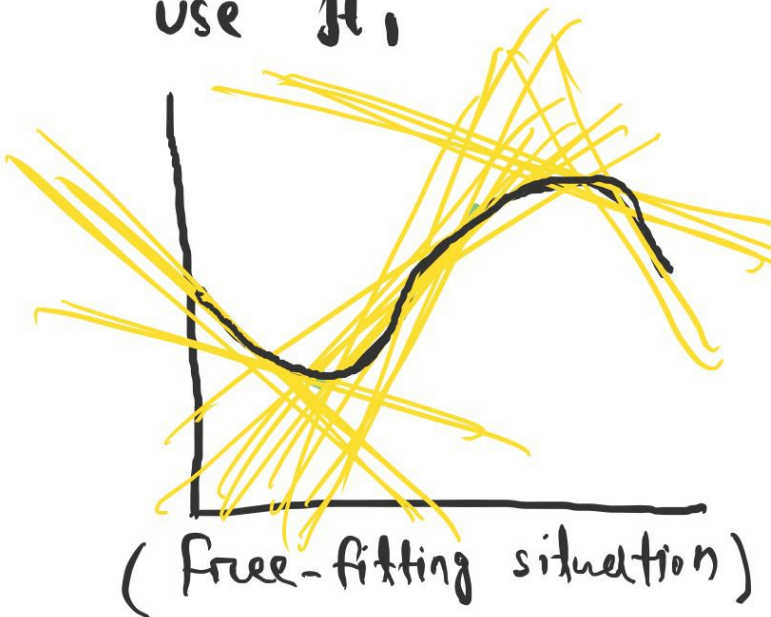
$$\text{minimize}_h E_{in}(h) + \Omega(h)$$

Some kind of measure on h . Ω is something that measure the complexity of h .

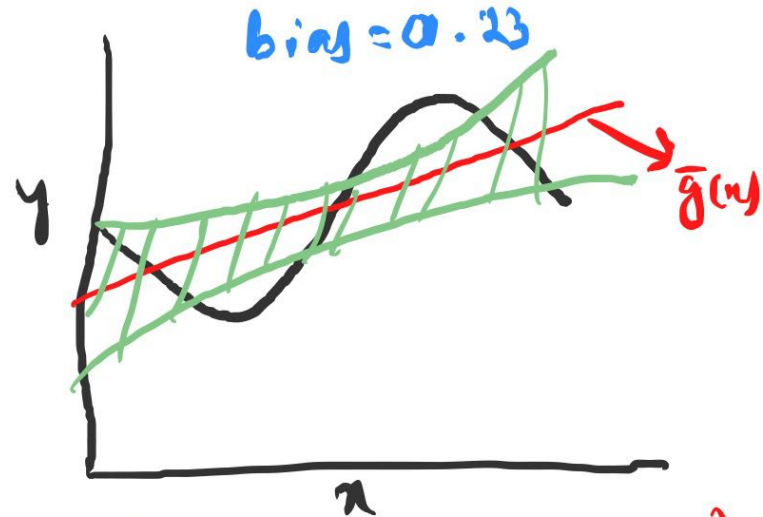
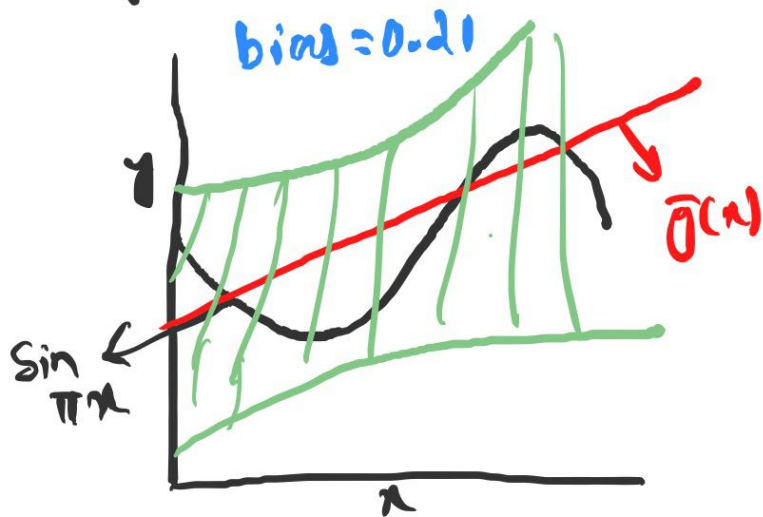
Motivation example :

Weight decay technique :

- measures the complexity of a hypothesis h by the size of the co-efficients used to represent h (e.g. linear model)
- Apply this to sine example, for $N=2$.
use \mathcal{H}_1



This weight decay technique reduces the variation by limiting the magnitude of your coefficients of the linear model.



(without regularization)

(with regularization)

Constant model $E_{out} = 0.75$

Unregularized model $E_{out} = 1.90$

Regularized model $E_{out} = 0.56$

→ Bias-Variance point of view: Improve variance but suffer from bias. overall is better. quality of the average fit

Q: Why need regularization?

→ The linear model is too sophisticated for the amount of data we have. Since a line can fit any 2 points.

→ The need of regularization depends on quantity and quality of data.

→ Given any two points, we can either choose

- a simple model e.g. constant model
- to constrain the model (e.g. weight decay method)

→ Constraining the model gives us more flexibility.

Regularization techniques :

1) Weight decay method

(constraint minimization)

2) Augmentation error

(unconstrained)

1) Weight decay method

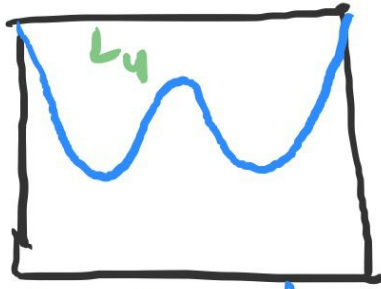
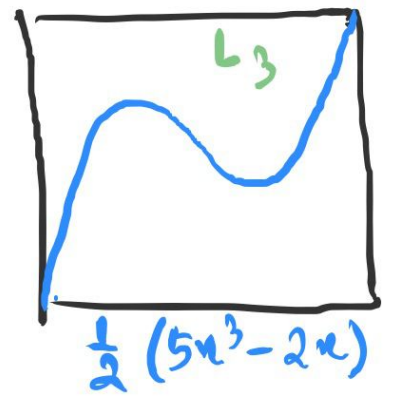
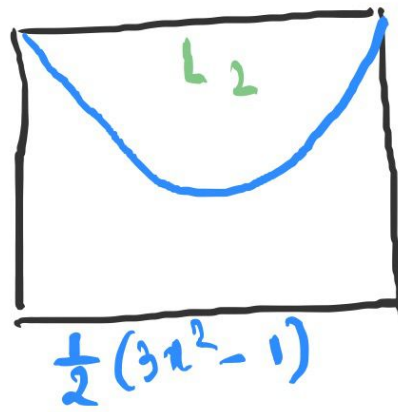
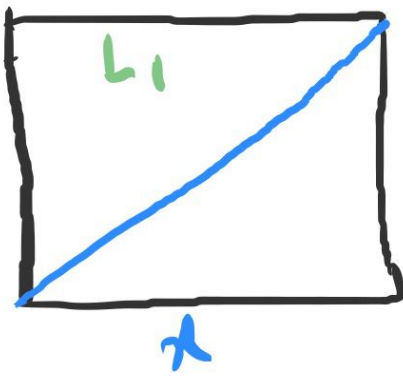
(soft order constraint)

→ Consider the following example

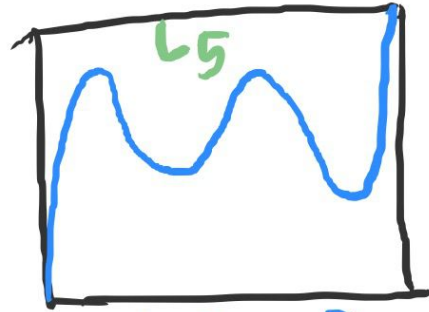
\mathcal{H} = set of polynomials in one variable

$x \in [-1, 1]$ e.g. $h(x) = 2x^2 + 3x + 7$

→ Express $h(x)$ using basis function. So the basis functions for polynomials are Legendre polynomials $L_q(x)$, $q = 1, 2, \dots$



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



$\frac{1}{8}(63x^5 + \dots)$

These polynomials are orthogonal.

So
$$h(x) = \sum_{q=1}^Q w_q L_q(x)$$

This model is linear.

If we define a non-linear transform Φ

$$z = \Phi(x) = \begin{bmatrix} L_1(x) \\ \vdots \\ L_Q(x) \end{bmatrix}$$

Then the hypothesis set is

$$\mathcal{H}_Q = \left\{ h \mid h(x) = w^T z = \sum_{q=0}^Q w_q L_q(x) \right\}.$$

Now we can define training error (for linear regression)

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

There are multiple ways of constraining the weights.

1) Hard constraint :

→ Force co-efficients to be zero.

→ e.g. $\mathcal{H}_2 = \left\{ w \mid w \in \mathcal{H}_{10} : \underline{w_q = 0}, \text{ for } \underline{q \geq 3} \right\}$

2) Soft constraint :

→ Force co-efficients to be small.

→ e.g. $\sum_{q=0}^Q w_q^2 \leq C$

→ This encourages weights to be small without changing the order of the polynomial by explicitly forcing some weights to 0.

VC perspective of soft order constraint:

The optimization problem is

$$\underset{w}{\text{minimize}} \quad E_{\text{in}}(w) \quad \text{Subject to} \quad w^T w \leq C.$$

We know $E_{\text{in}}(w) = \frac{1}{N} \|Zw - y\|_2^2$

\downarrow \rightarrow labels

transformed input
vector put into a matrix

The soft order constrained hypothesis is

$$\mathcal{H}(C) = \{ h \mid h(x) = w^T z, w^T w \leq C \}$$

The optimization is equivalent to saying
minimize

$$E_{\text{in}} \text{ over } \mathcal{H}(C).$$

\rightarrow If $C_1 < C_2$, then $\mathcal{H}(C_1) \subset \mathcal{H}(C_2)$.

and $d_{\text{VC}}(\mathcal{H}(C_1)) \leq d_{\text{VC}}(\mathcal{H}(C_2))$

\rightarrow We should expect better generalization
with $\mathcal{H}(C_1)$.

Augmented error

(Unconstrained minimization)

$$E_{\text{aug}}(w) = E_{\text{in}}(w) + \lambda w^T w$$

$\lambda > 0$ is a free parameter.

- Unconstrained minimization is often easier than constrained minimization.
- For a given C , soft order constraint corresponds to selecting a hypothesis from a smaller set $\mathcal{H}(C)$.
- For augmented error, you need to find the optimal parameter λ^* .

VC perspective of augmented error :

The augmented error for a hypothesis $h \in \mathcal{H}$ is

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

Here, $\Omega(h) = w^T w$.

There are two components of the penalty:

- 1) The regularizer $\Omega(h)$ which penalizes a particular property of h .
- 2) Regularization parameter λ (controls the amount of regularization)