

Foundations of Machine Learning

AI2000 and AI5000

FoML-09

Underfitting and Overfitting
Regularized Least Squares

Dr. Konda Reddy Mopuri

Department of AI, IIT Hyderabad
July-Nov 2025



భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad



So far in FoML

- What is ML and the learning paradigms
- Probability refresher
- MLE, MAP, and fully Bayesian treatment
- Linear Regression with basis functions - geometric interpretation



Under and Over fitting



భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad



Linear Regression

- Complex functions can be fit to the data
 - Using the basis functions



భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad



Linear Regression with basis functions

- Complex functions can be fit to the data
 - Using the basis functions
- There are some choices (hyper parameters) to be made
 - What kind of basis functions?
 - How many of them?



Linear Regression with basis functions

- Complex functions can be fit to the data
 - Using the basis functions
- There are some choices (hyper parameters) to be made
 - What kind of basis functions?
 - How many of them?
- They have consequences
 - over/under fitting

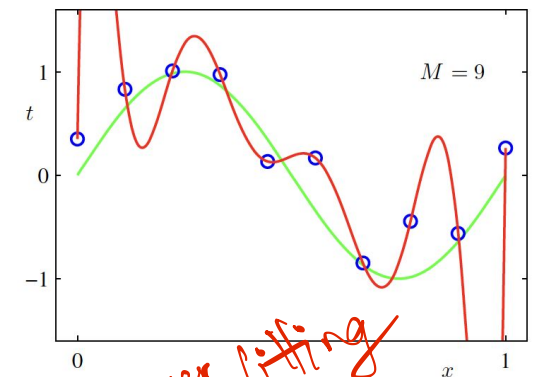
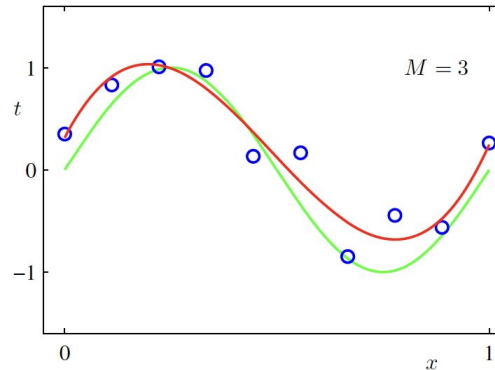
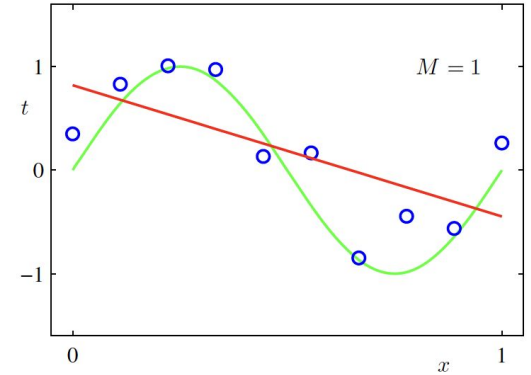
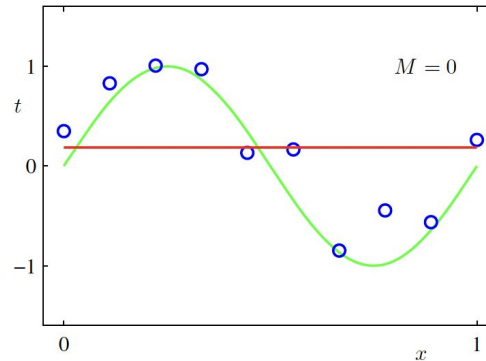


Polynomial basis functions

Target $t = \sin(2\pi x) + \sigma \cdot \epsilon$

Noise $\epsilon \sim N(0, 1)$

$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{j=0}^M w_jx^j$$



underfitting

overfitting



భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad

DiL
Data-driven Intelligence
& Learning Lab

How to spot Under/Over fitting?



భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad



How to spot Under/Over fitting?

- Can the weight values provide some insight?



భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad



How to spot Under/Over fitting?

- Can the weight values provide some insight?

	$M = 0$	$M = 1$	$M = 3$	$M = 9$
w_0^*	0.19	0.82	0.31	0.35
w_1^*		-1.27	7.99	232.37
w_2^*			-25.43	-5321.83
w_3^*			17.37	48568.31
w_4^*				-231639.30
w_5^*				640042.26
w_6^*				-1061800.52
w_7^*				1042400.18
w_8^*				-557682.99
w_9^*				125201.43

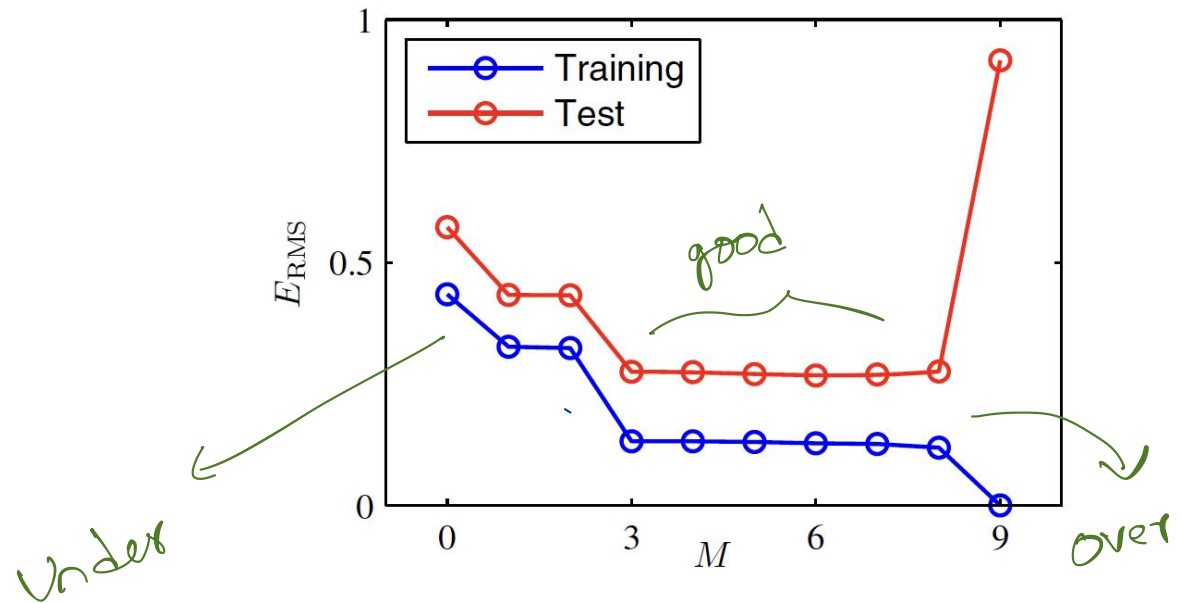
very Big



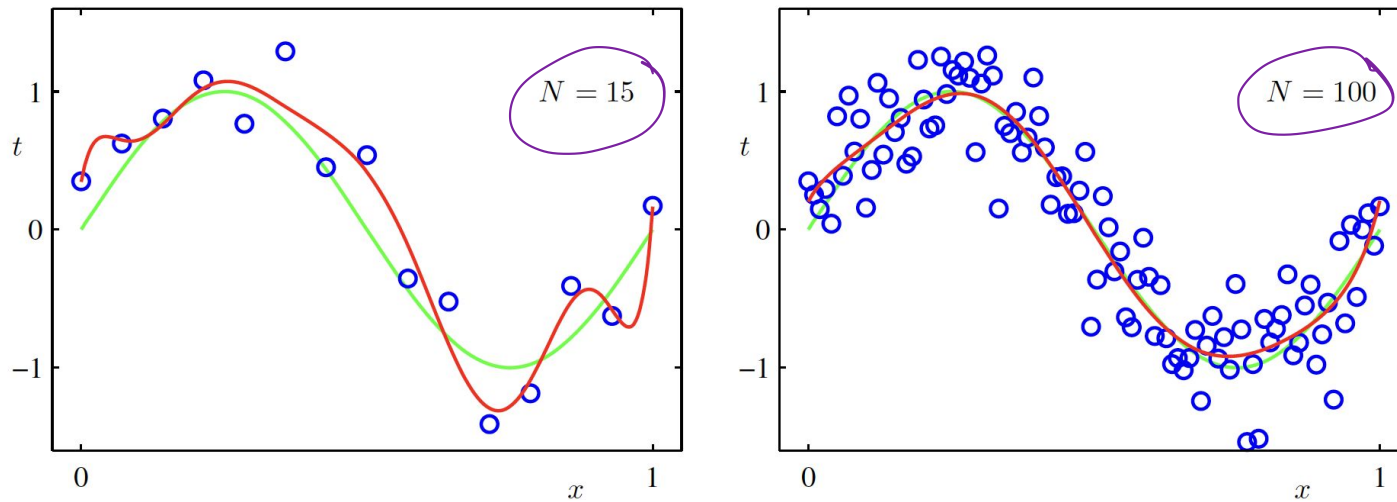
Why didn't our $M = 9$ model realize a mapping similar to that of $M = 3$ model?

Better way to spot under/over fitting

- Plot the error (E_{RMS})



Effect of dataset size on overfitting



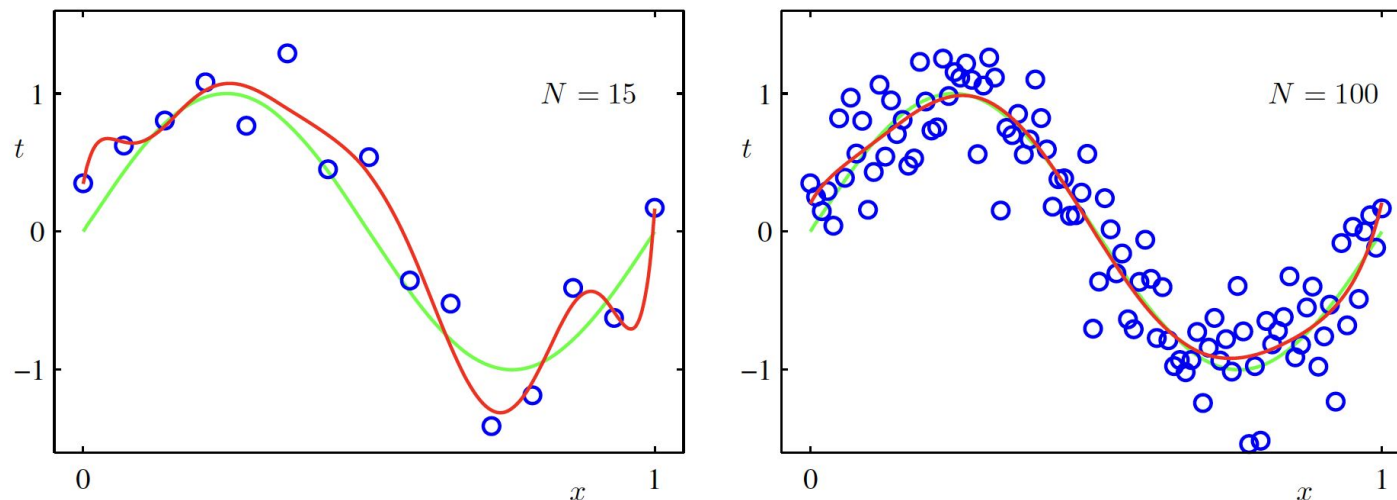
M = 9 model fitting to dataset of different sizes

$$N \downarrow \begin{bmatrix} | & | & | & | & | & | & | & | & | & | \end{bmatrix} \Phi \quad \underline{w} = \hat{y} \approx y$$

$N \times M$



Effect of dataset size on overfitting



$M = 9$ model fitting to dataset of different sizes

For a given model complexity, the overfitting problem becomes less severe as the dataset size increases



What if it is not easy to collect a lot of data?

- More data helps to avoid overfitting
- But, it may be challenging to collect a lot of data
- Then, what?

Regularization



భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad



How to spot Under/Over fitting?

- Large weight values indicate overfitting
- Higher complexity models lead to overfitting

	$M = 0$	$M = 1$	$M = 3$	$M = 9$
w_0^*	0.19	0.82	0.31	0.35
w_1^*		-1.27	7.99	232.37
w_2^*			-25.43	-5321.83
w_3^*			17.37	48568.31
w_4^*				-231639.30
w_5^*				640042.26
w_6^*				-1061800.52
w_7^*				1042400.18
w_8^*				-557682.99
w_9^*				125201.43



Regularized least squares

- In case of smaller datasets, instead of manually restricting the number of parameters
 - Add a penalty to avoid large weight values → 'weight decay' regularization

called



Regularized least squares

- In case of smaller datasets, instead of manually restricting the number of parameters
 - Add a penalty to avoid large weight values → 'weight decay' regularization

$$\tilde{E}(\mathbf{w}) = \underbrace{\frac{1}{2} \sum_{i=1}^N \{t_i - \mathbf{w}^T \phi(\mathbf{x}_i)\}^2}_{\text{MSE}} + \underbrace{\lambda \sum_{i=1}^M w_i^2}_{\text{regularization weight-decay or parameter shrinkage}}$$

new objective



Regularized least squares

- In case of smaller datasets, instead of manually restricting the number of parameters
 - Add a penalty to avoid large weight values

Ridge
Regression

$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N \{t_i - \mathbf{w}^T \phi(\mathbf{x}_i)\}^2 + \frac{\lambda}{2} \sum_{i=1}^{M-1} w_i^2$$

Bias term w_0 may not be included in the regularization



Regularized least squares

$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N \{t_i - \mathbf{w}^T \phi(\mathbf{x}_i)\}^2 + \frac{\lambda}{2} \sum_{i=1}^{M-1} \mathbf{w}_i^2$$

- Looks similar to what we saw during the MAP discussion

MAP objective = $\frac{\beta}{2} \sum_{i=1}^N \left(\text{residual} \right)^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$

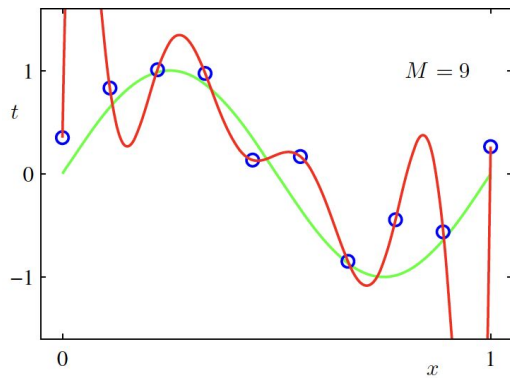
$\lambda = \alpha / \beta$

$= \beta \left[\text{MSE} + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \right]$

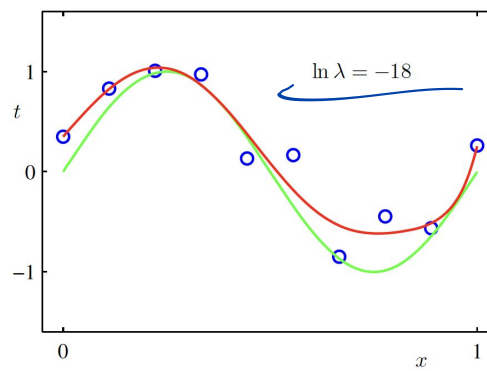
$= \beta \tilde{E}(\mathbf{w})$



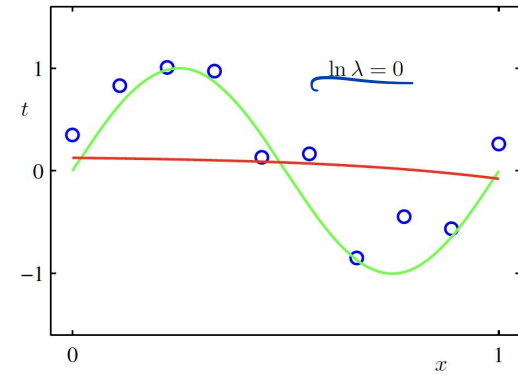
Regularized Least Squares



w/o regularization



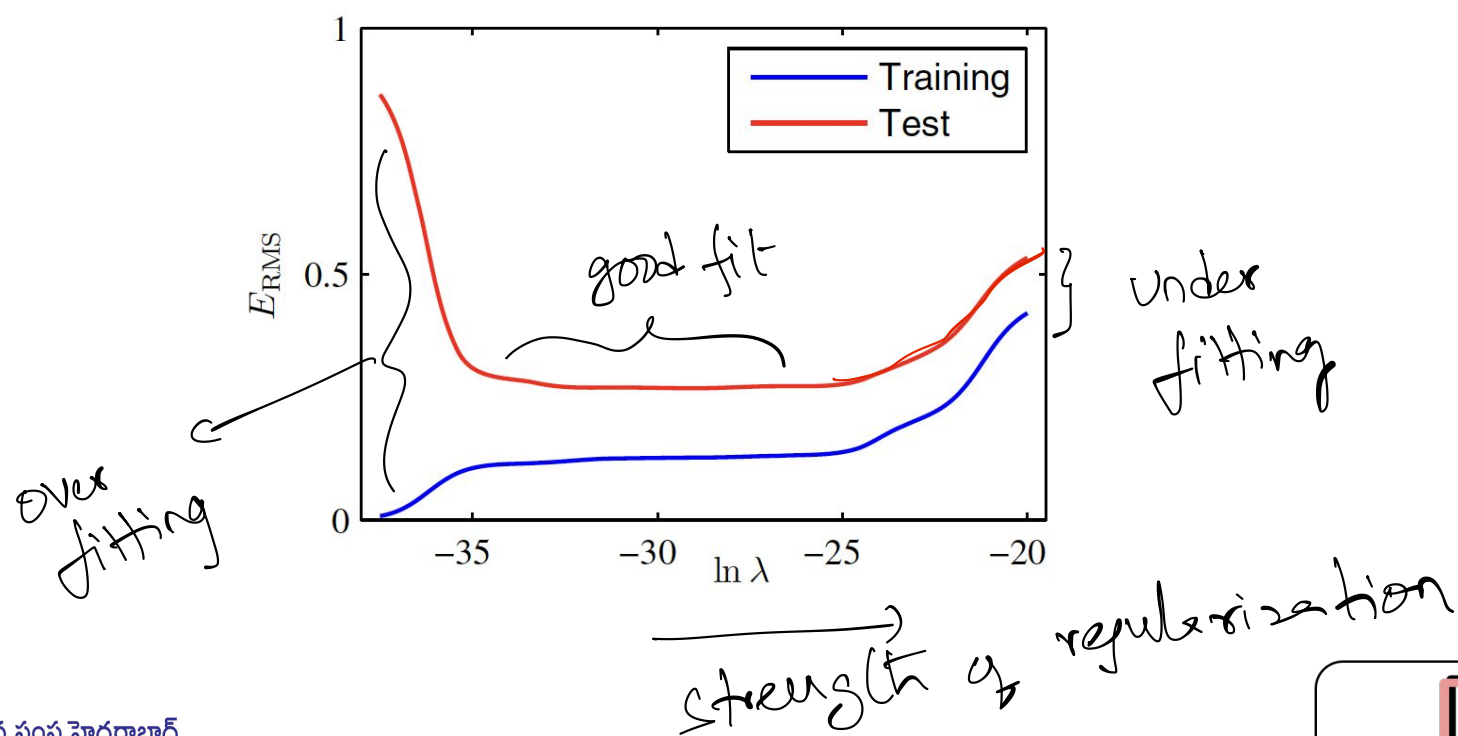
moderate
regularization



Heavy
regularization



Regularized Least Squares



More general form of the regularization

$$\tilde{E}(\mathbf{w}) = \underbrace{\frac{1}{2} \sum_{i=1}^N \{t_i - \mathbf{w}^T \phi(\mathbf{x}_i)\}^2}_{\text{Empirical Error}} + \underbrace{\frac{\lambda}{2} \sum_{i=1}^{M-1} |\mathbf{w}_i|^q}_{\text{Regularization}}$$

- When $q = 2 \rightarrow l_2$ norm penalty on the parameters



More general form of the regularization

$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N \{t_i - \mathbf{w}^T \phi(\mathbf{x}_i)\}^2 + \frac{\lambda}{2} \sum_{i=1}^{M-1} |\mathbf{w}_i|^q$$

- When $q = 2 \rightarrow l_2$ norm ^{square} penalty on the parameters
- When $q = 1 \rightarrow l_1$ norm penalty (also, called Lasso)



Geometric interpretation

Equivalent to minimizing

$$\frac{1}{2} \sum_{i=1}^N \{t_i - \mathbf{w}^T \phi(\mathbf{x}_i)\}^2 \quad \text{with} \quad \sum_{j=1}^M |w_j|^q \leq \eta$$

objective *feasible set*

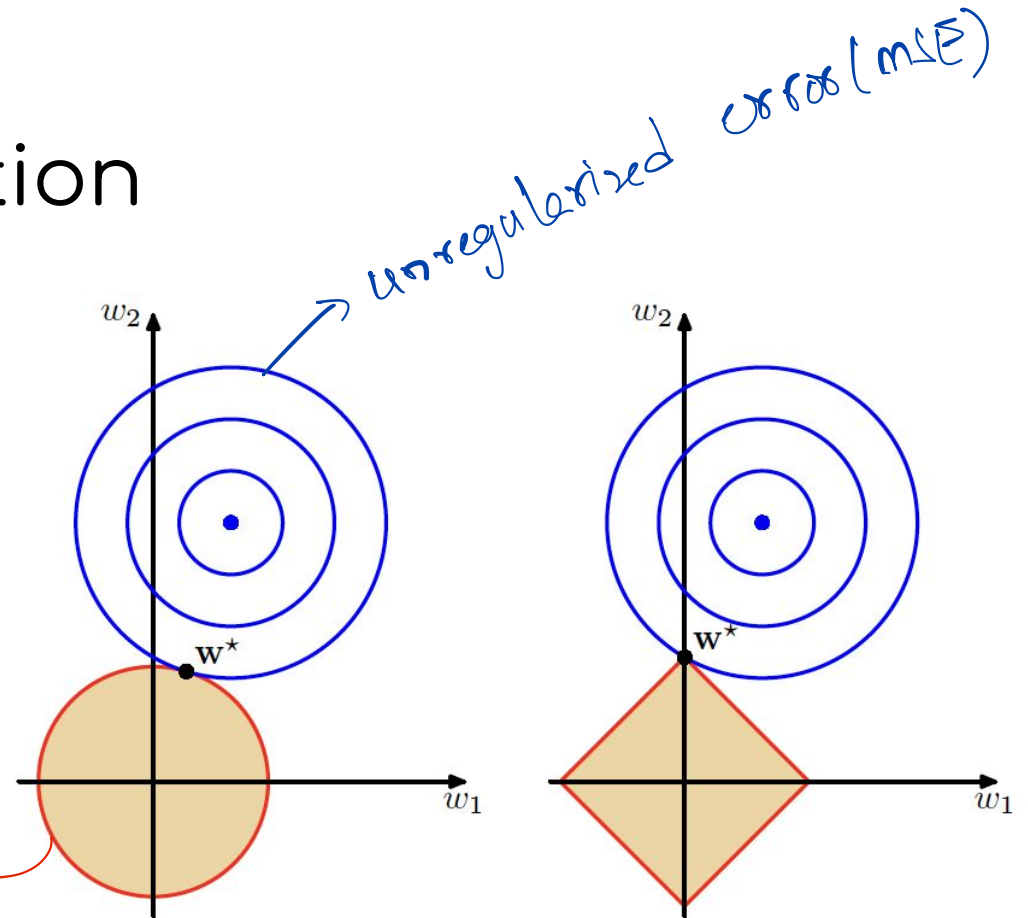


Geometric interpretation

Equivalent to minimizing

$$\frac{1}{2} \sum_{i=1}^N \{t_i - \mathbf{w}^T \phi(\mathbf{x}_i)\}^2 \quad \text{with}$$

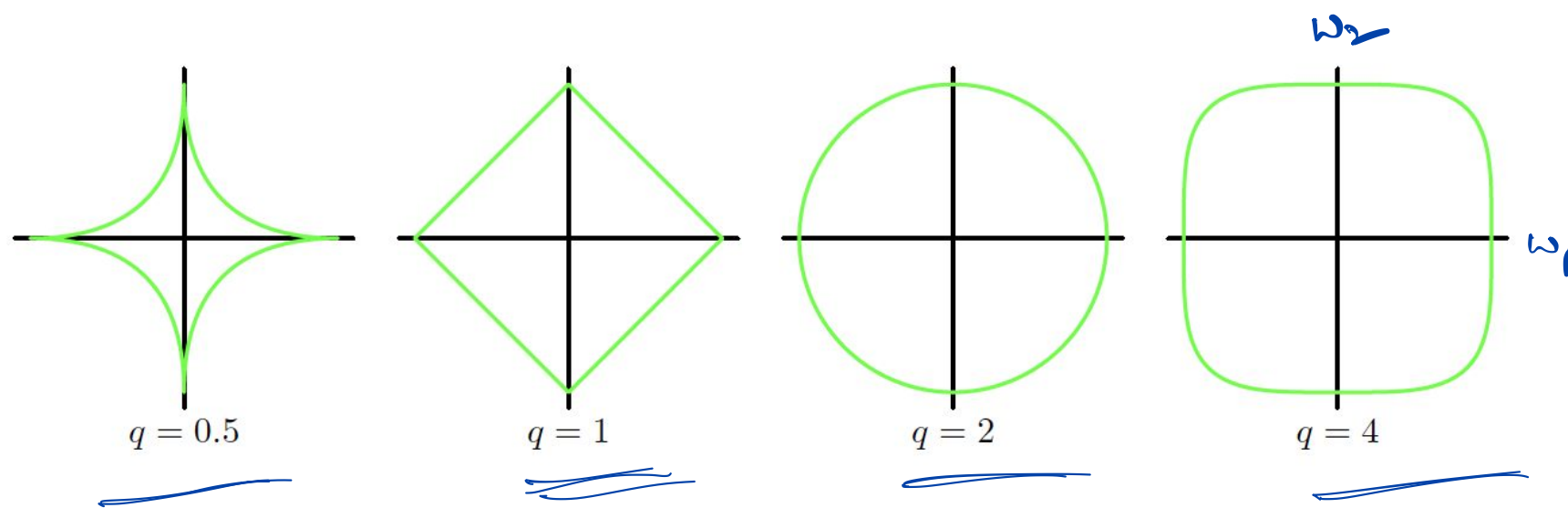
$$\sum_{j=1}^M |w_j|^q \leq \eta$$



Plots of the 'unregularized' error and the 'regularization' term for $q=1$ and 2



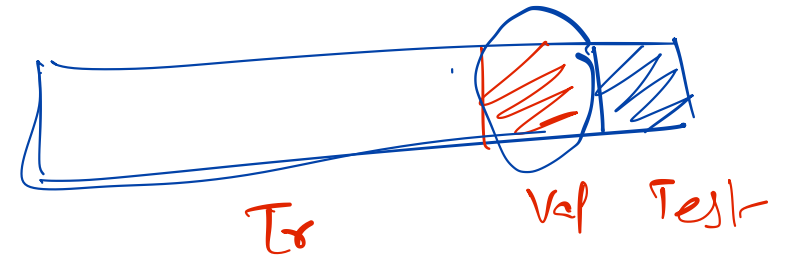
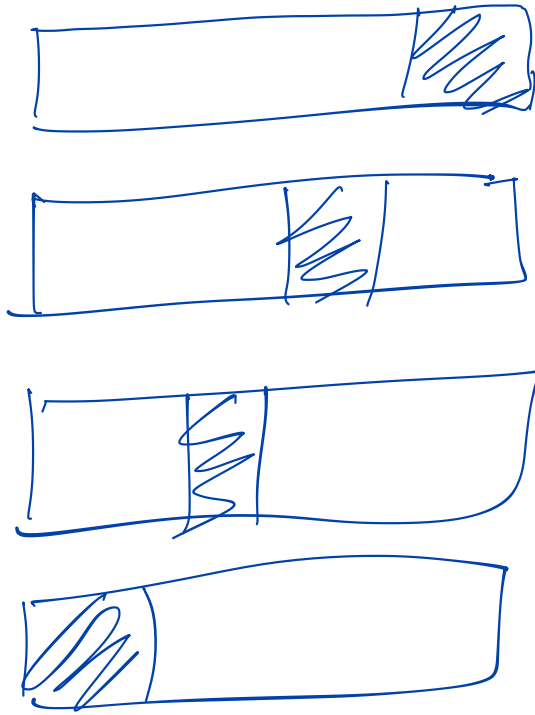
Geometric interpretation



Contours of regularization term for different values of q



Rough work



Model Selection via
Cross validation



భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad



Next Model selection



భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad

