



# CPE/EE 695: Applied Machine Learning

## ***Lecture 2 -1 : Polynomial Fitting***

Dr. Shucheng Yu, Associate Professor  
Department of Electrical and Computer Engineering  
Stevens Institute of Technology



# Elements for the learning problems

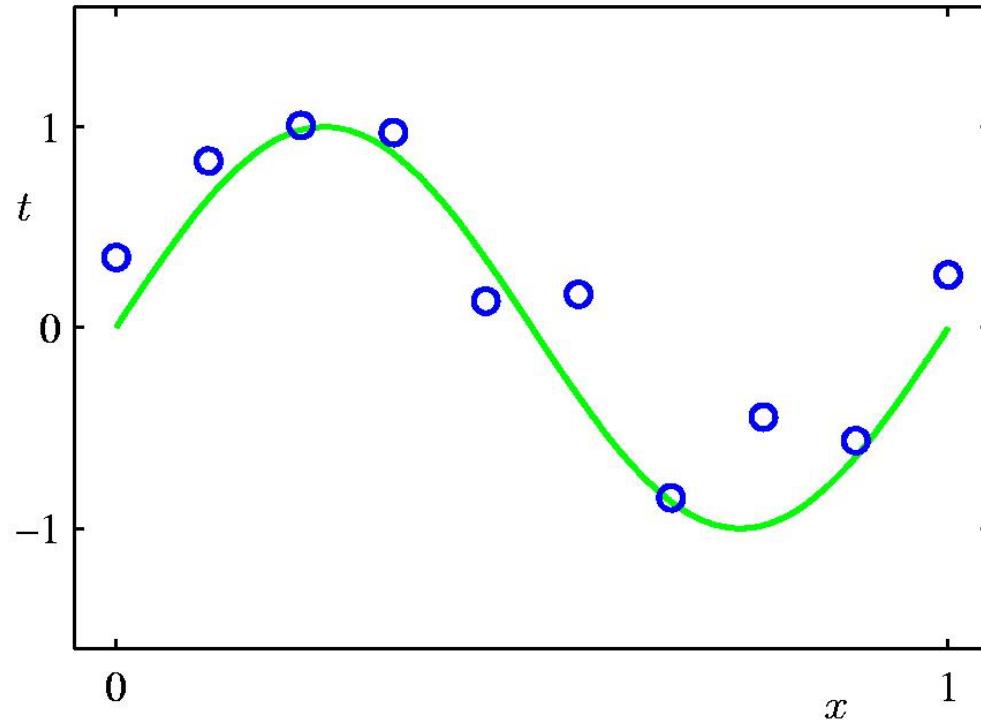
Learning = Improving with experience at some task

- Improve over task T,
- with respect to performance measure P,
- based on experience E.

E.g., Learn to play checkers

- T: Play checkers
- P: % of games won in world tournament
- E: opportunity to play against self

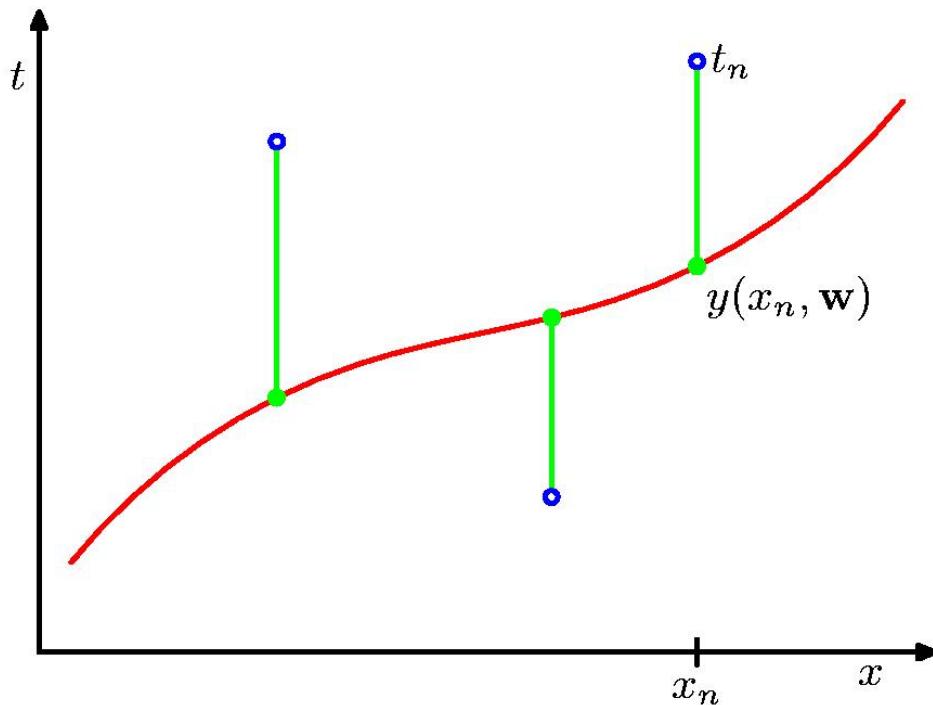
# Polynomial Curve Fitting



$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \dots + w_M x^M = \sum_{j=0}^M w_j x^j$$

How to measure the fit between model and training data?

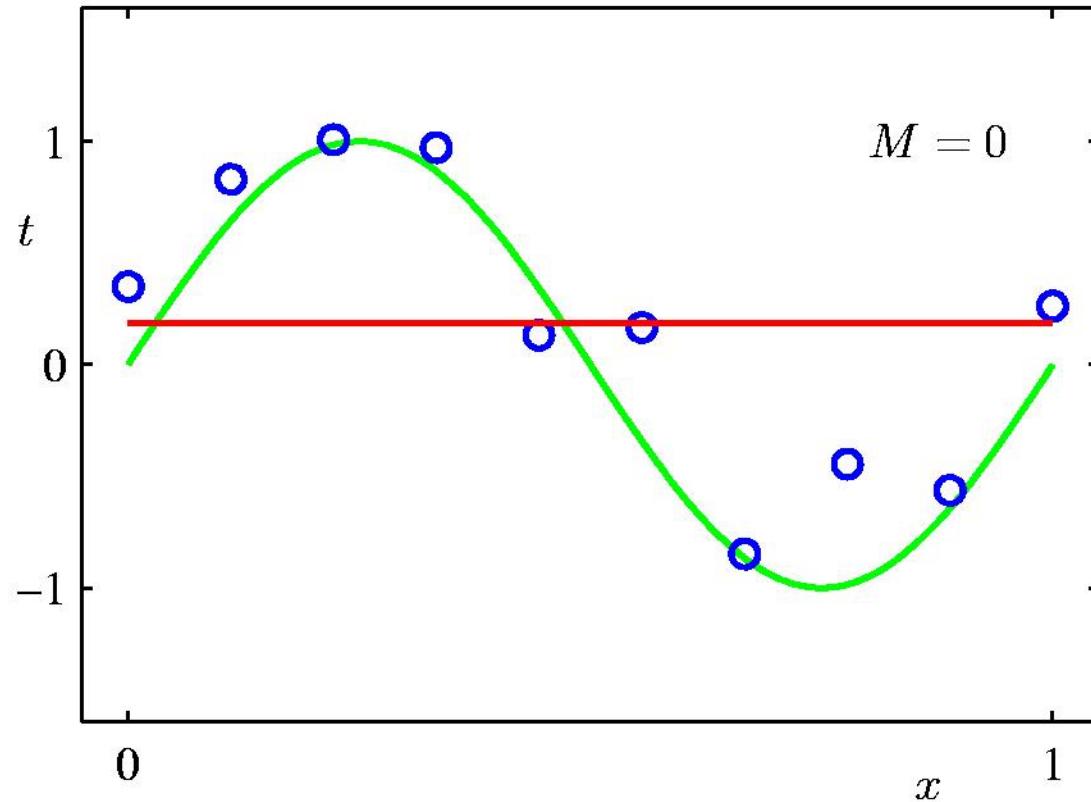
# Sum-of-Squares Error Function



$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2$$

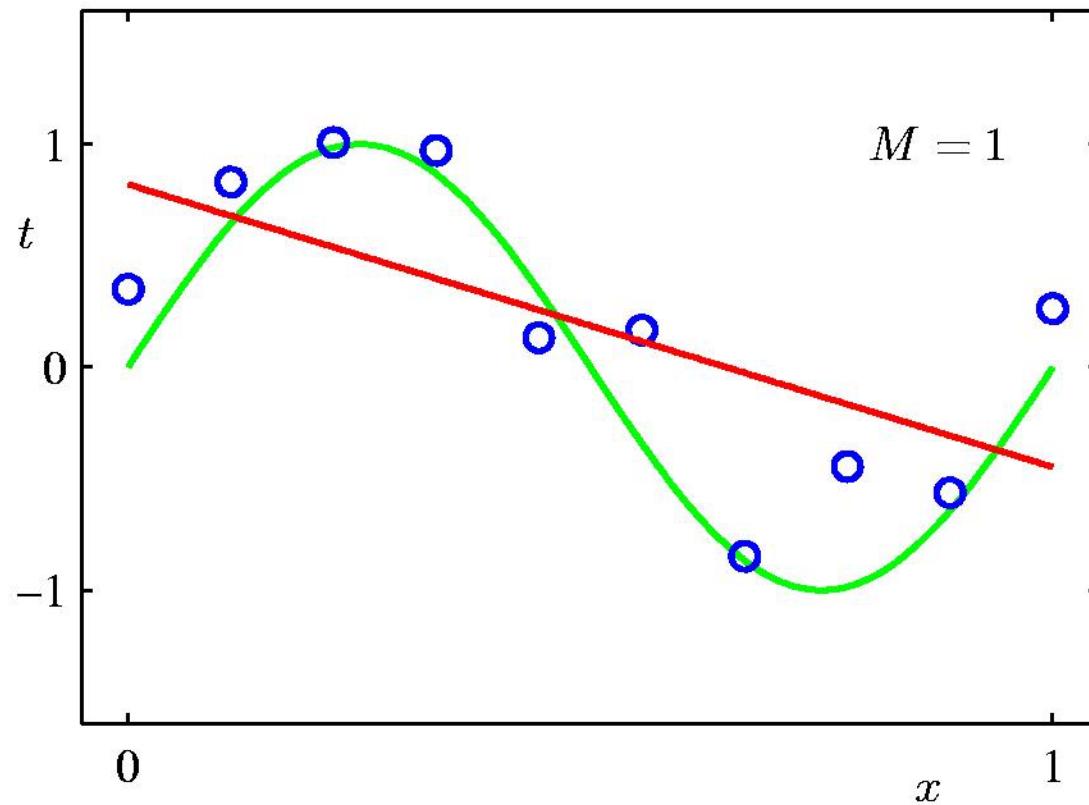
Minimize  $E(\mathbf{w})$  for unknown  $\mathbf{w}$ . (maximum likelihood)

# 0<sup>th</sup> Order Polynomial

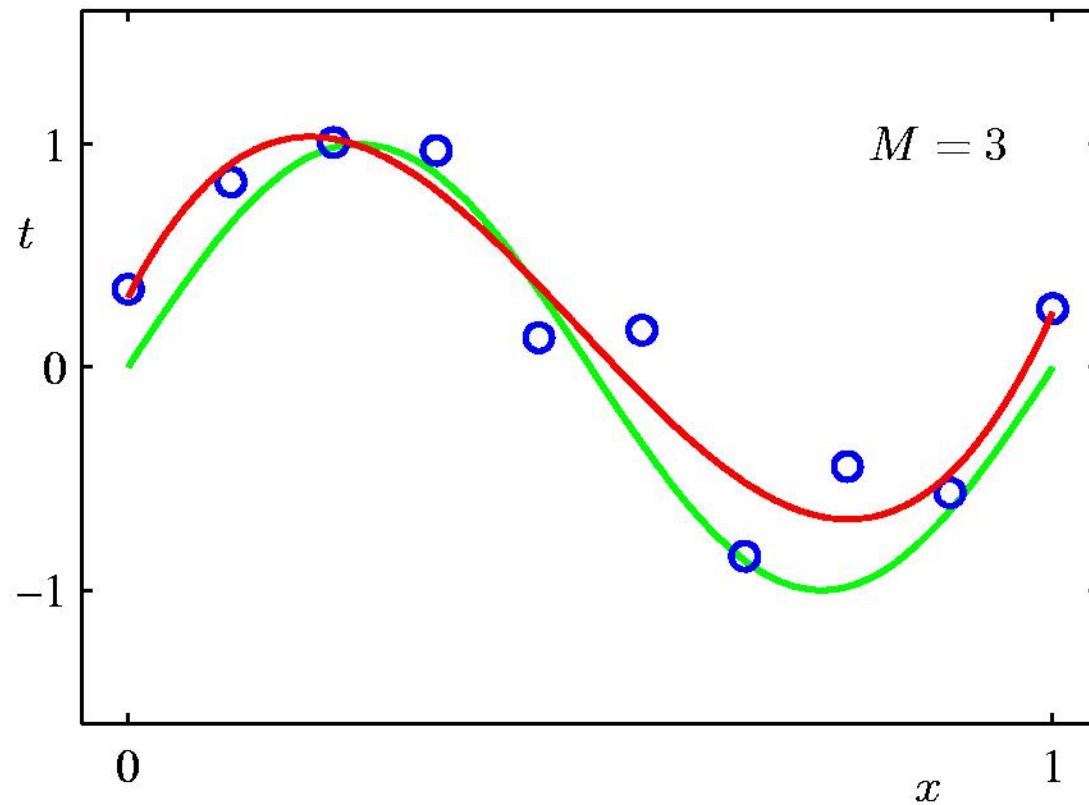


Model selection: how to choose the order M?

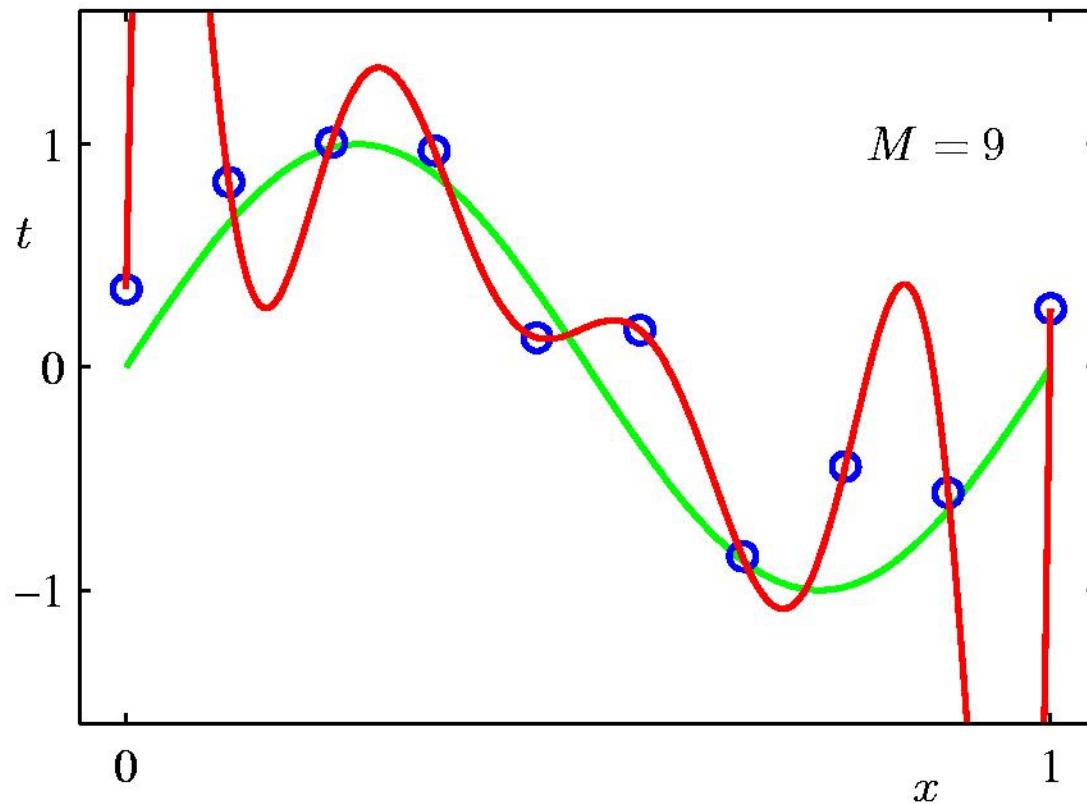
# 1<sup>st</sup> Order Polynomial



# 3<sup>rd</sup> Order Polynomial

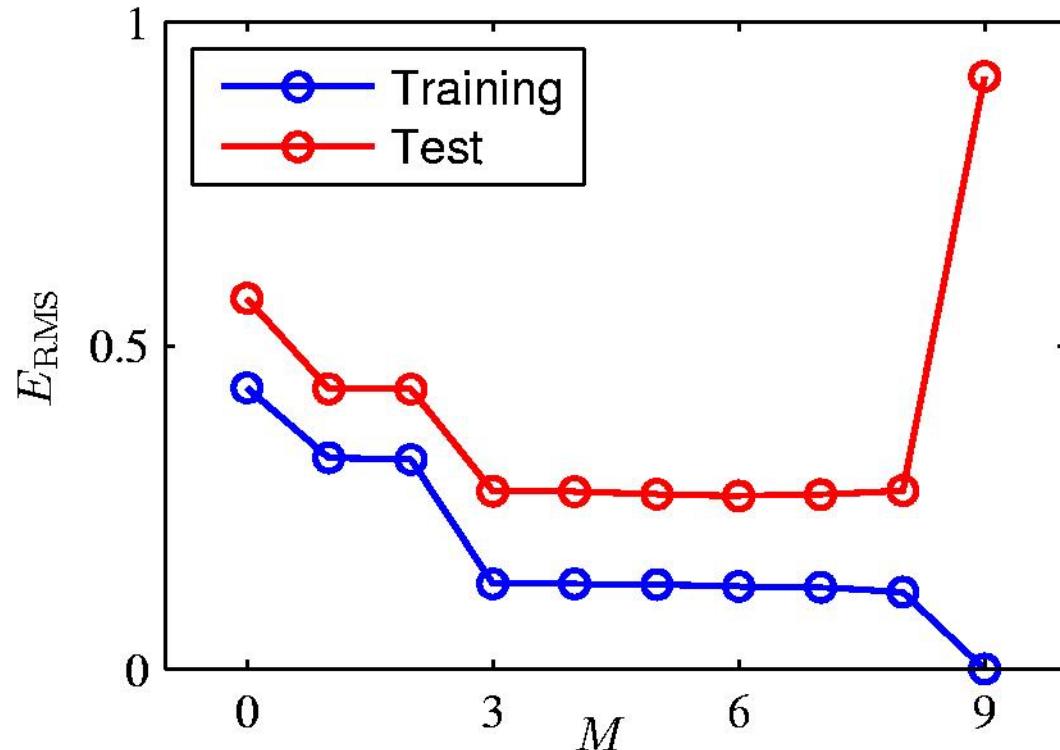


# 9<sup>th</sup> Order Polynomial



M=9: perfectly fit for training data set. Question: the larger M the better?

# Over-fitting



Root-Mean-Square (RMS) Error:  $E_{\text{RMS}} = \sqrt{2E(\mathbf{w}^*)/N}$

$M=9$ : good for training data, not for test data. What is under-fitting?  
**Bias-Variance Problem**



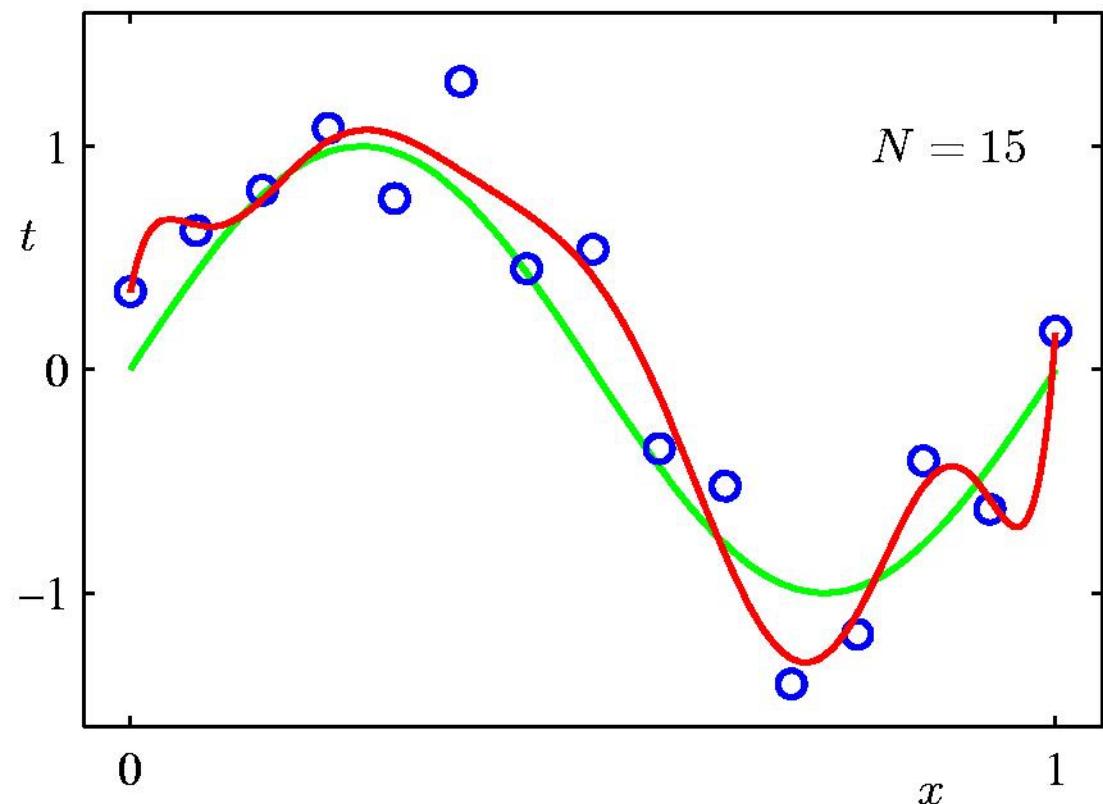
# Polynomial Coefficients

	$M = 0$	$M = 1$	$M = 3$	$M = 9$
$w_0^*$	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

coefficients increases as M getting larger (larger oscillations).

Data Set Size:  $N = 15$

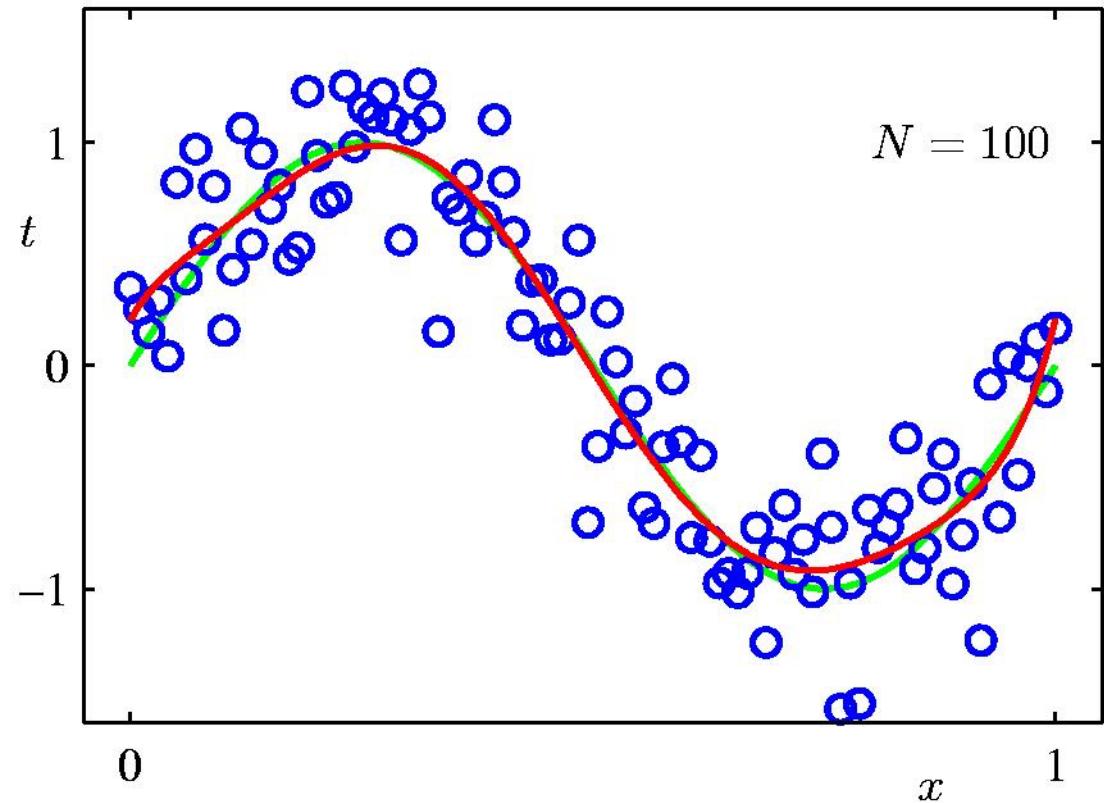
9<sup>th</sup> Order Polynomial



Overfitting less severe as data set size increases

Data Set Size:  $N = 100$

9<sup>th</sup> Order Polynomial



The **larger** the data set, the more **complex** model we can afford



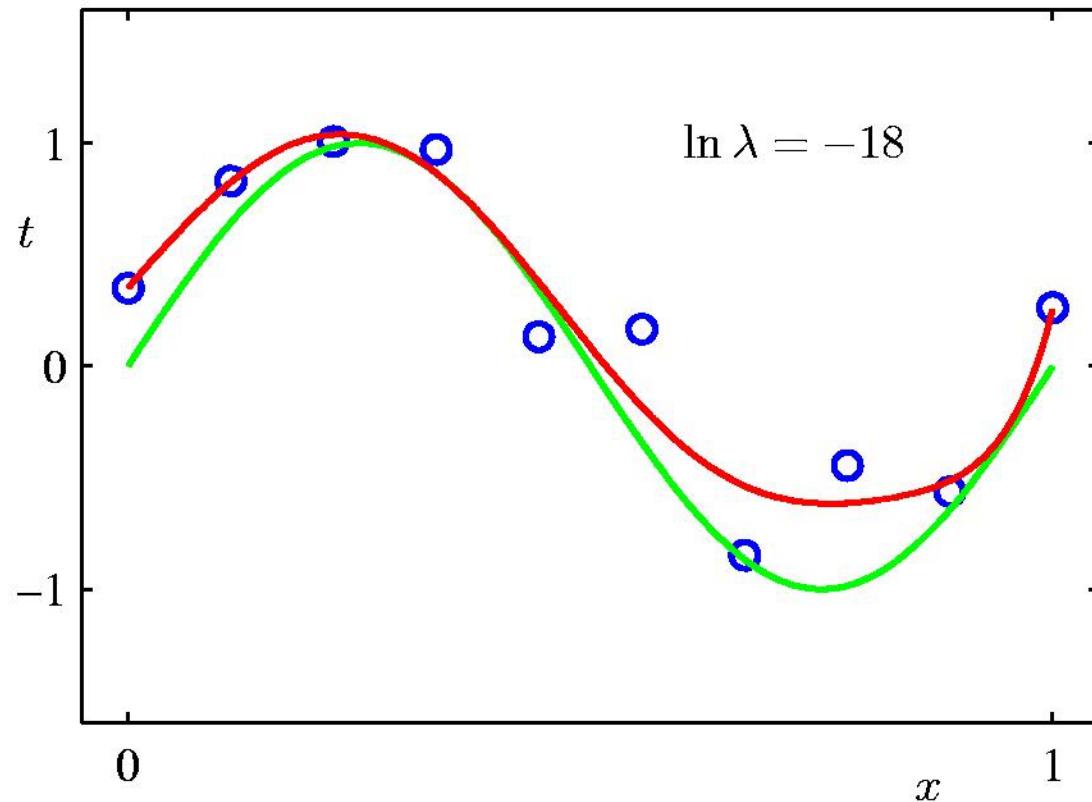
# Regularization

- Penalize large coefficient values

$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

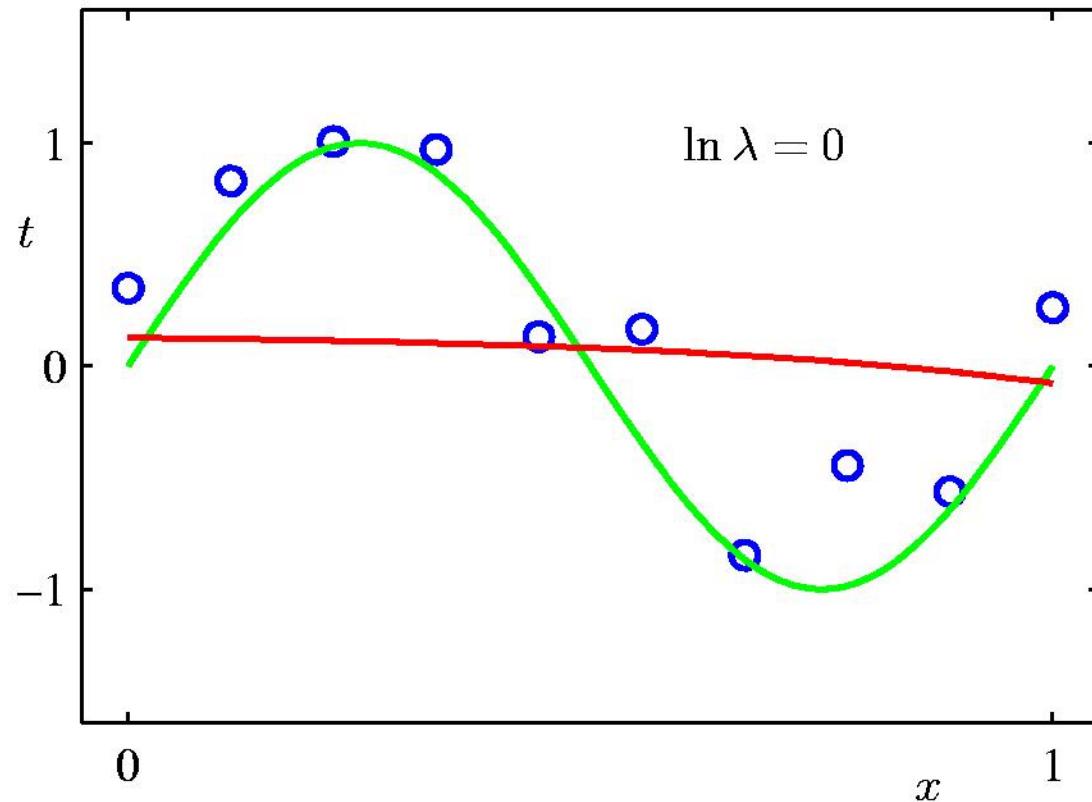
To support complex model under limited data size for maximum likelihood approach.

**Regularization:**  $\ln \lambda = -18$

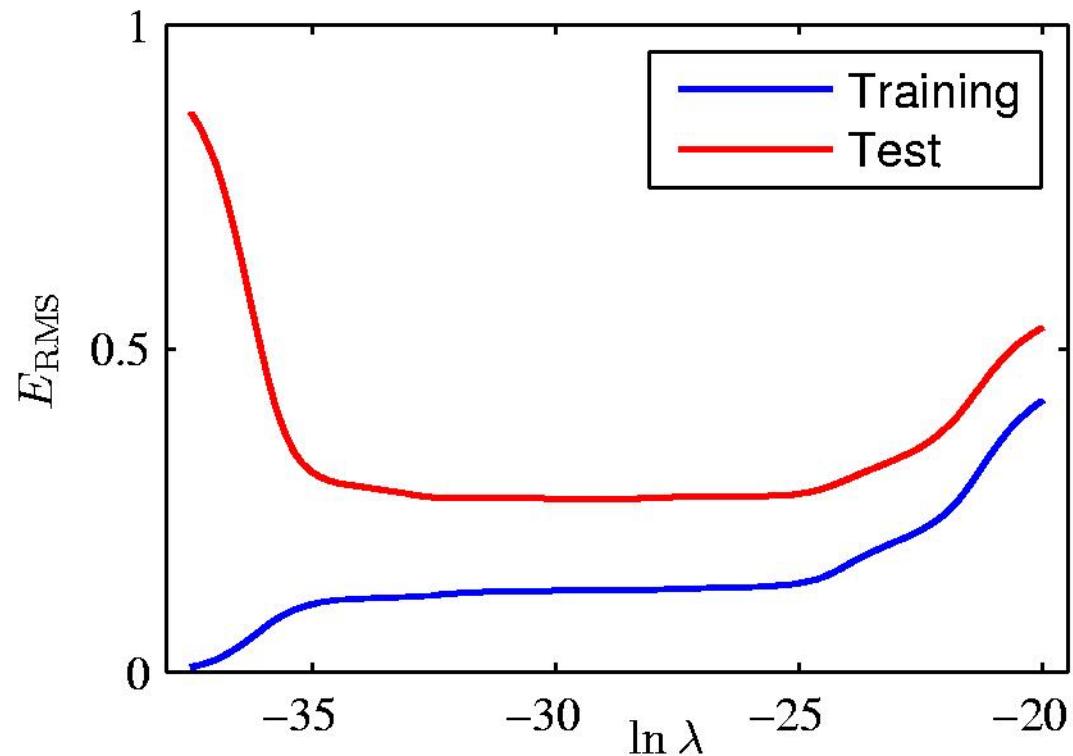


$M = 9$

# Regularization: $\ln \lambda = 0$



# Regularization: $E_{\text{RMS}}$ vs. $\ln \lambda$





# Polynomial Coefficients

	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
$w_0^*$	0.35	0.35	0.13
$w_1^*$	232.37	4.74	-0.05
$w_2^*$	-5321.83	-0.77	-0.06
$w_3^*$	48568.31	-31.97	-0.05
$w_4^*$	-231639.30	-3.89	-0.03
$w_5^*$	640042.26	55.28	-0.02
$w_6^*$	-1061800.52	41.32	-0.01
$w_7^*$	1042400.18	-45.95	-0.00
$w_8^*$	-557682.99	-91.53	0.00
$w_9^*$	125201.43	72.68	0.01



# Acknowledgement

Significant part of the slides were borrowed from Dr. Christopher M. Bishop's material at  
<https://www.microsoft.com/en-us/research/people/cmbishop/#!prml-book?from=http%3A%2F%2Fresearch.microsoft.com%2F%7Ecmbishop%2Fprml>;

For details please read Chapter 1 of textbook "Pattern Recognition and Machine Learning", by Christopher Bishop.



**STEVENS**  
INSTITUTE *of* TECHNOLOGY  
THE INNOVATION UNIVERSITY®

**stevens.edu**

