

Artificial Intelligence

CE-417, Group 1

Computer Eng. Department

Sharif University of Technology

Spring 2024

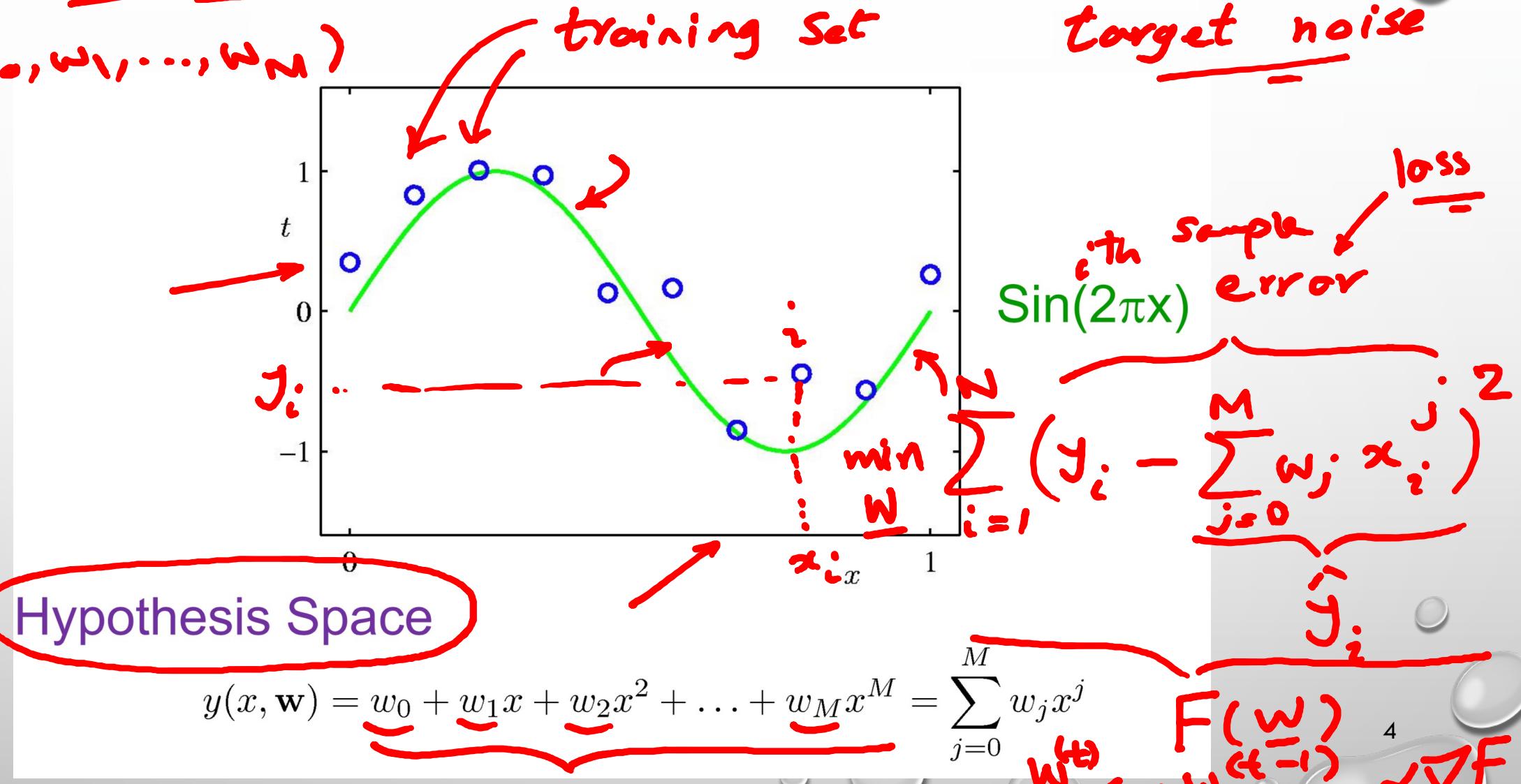
By Mohammad Hossein Rohban, Ph.D.

Courtesy: Most slides are adopted from CSE-573 (Washington U.), original
slides for the textbook, and CS-188 (UC. Berkeley).

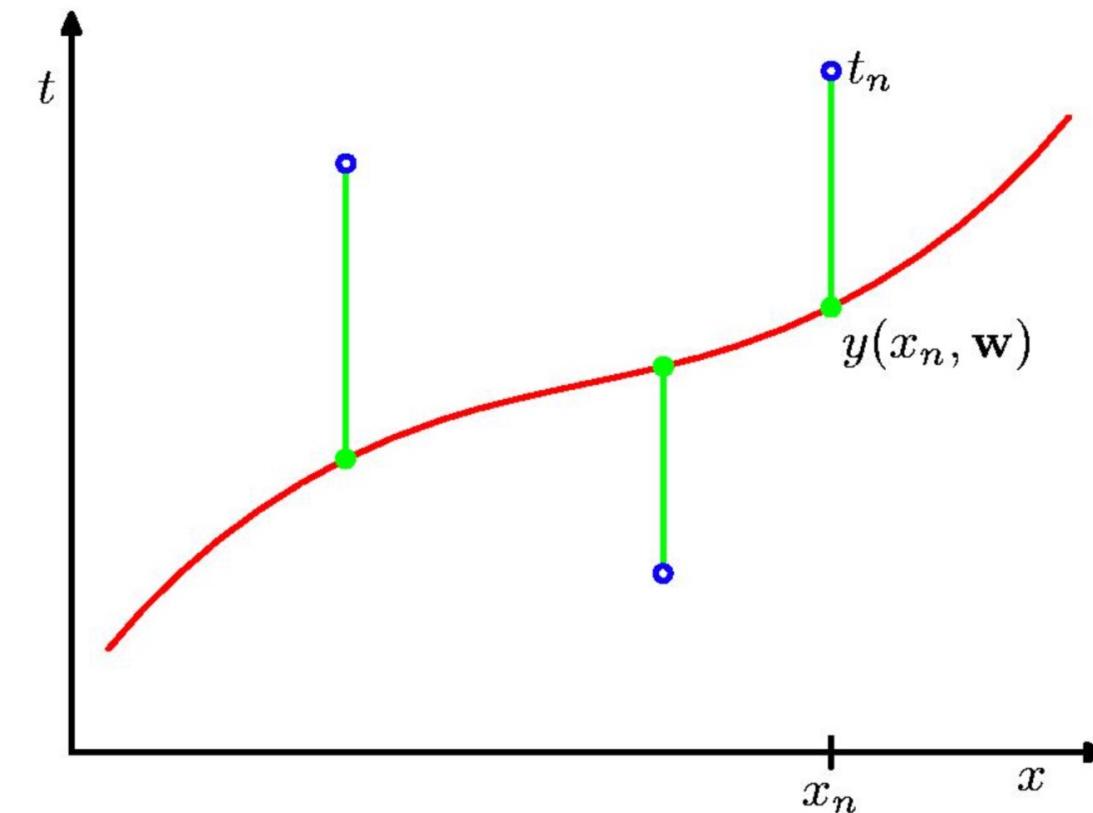
Regression

Inductive Bias Polynomial Curve Fitting

$$\underline{w} = (w_0, w_1, \dots, w_M)$$



Sum-of-Squares Error Function



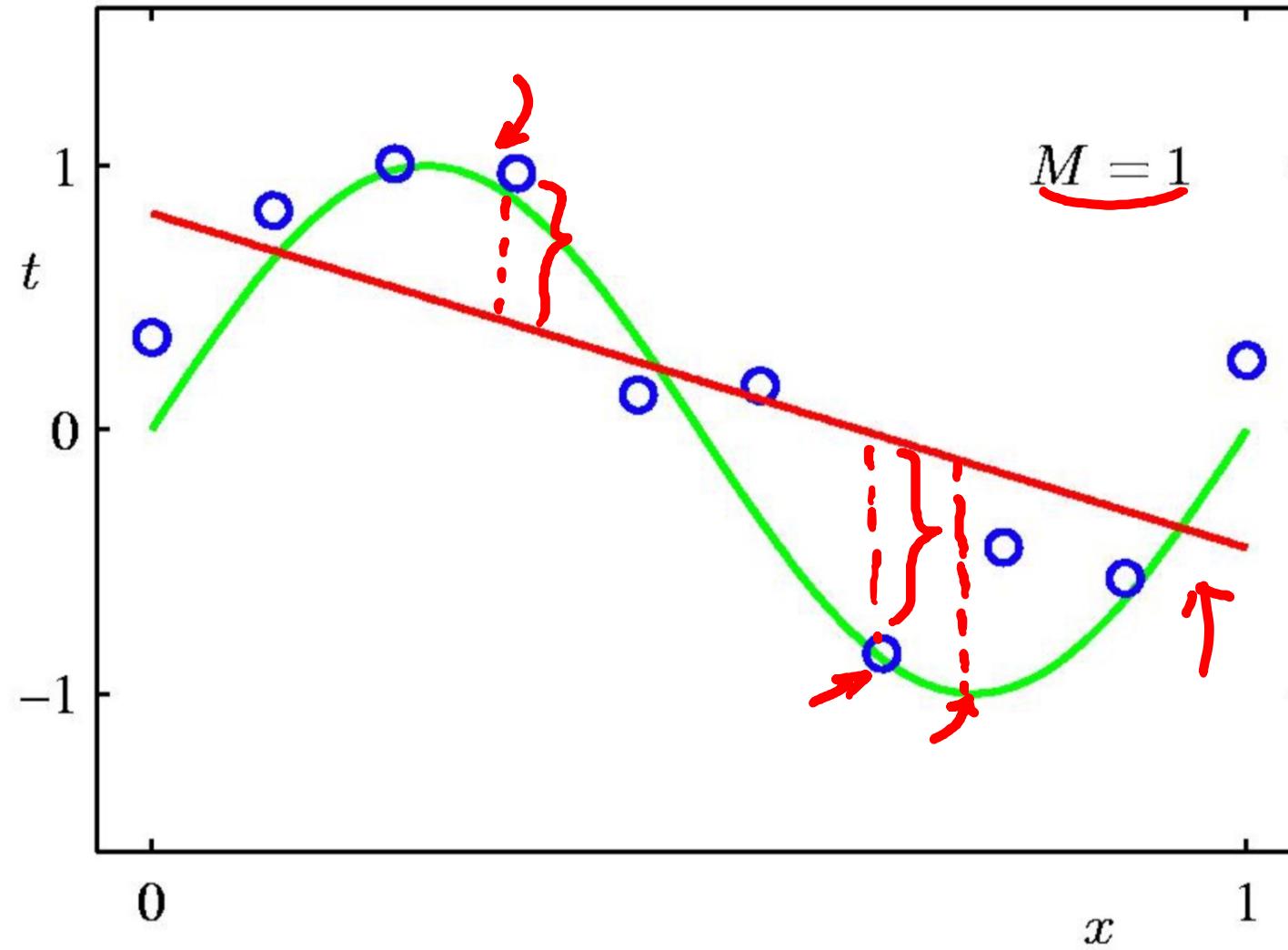
$$\underline{E(\mathbf{w})} = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2$$

1st Order Polynomial

Underfit

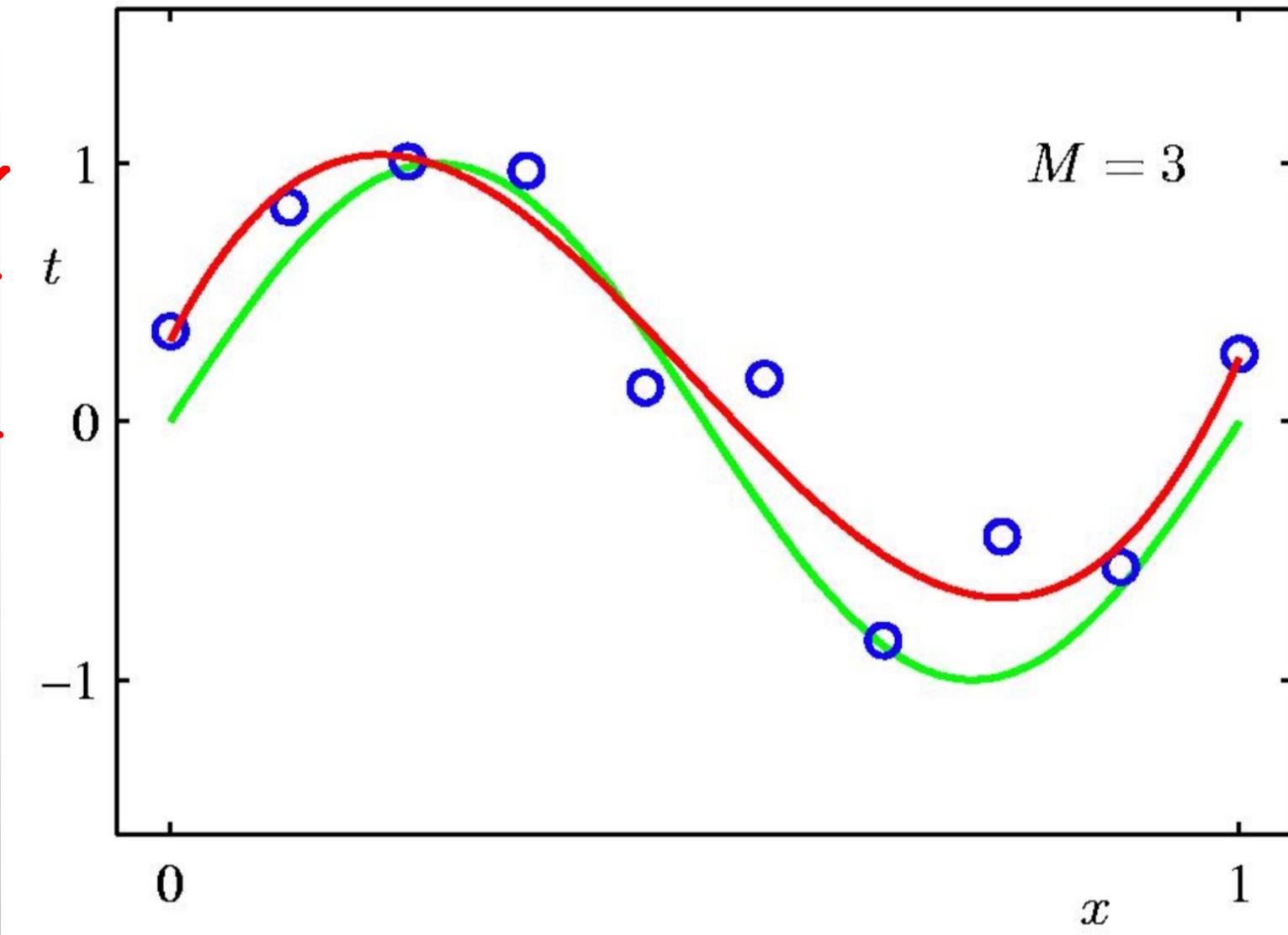
trn err \uparrow

tst err \uparrow



3rd Order Polynomial

trn err ↓
tst err ↓
Right fit

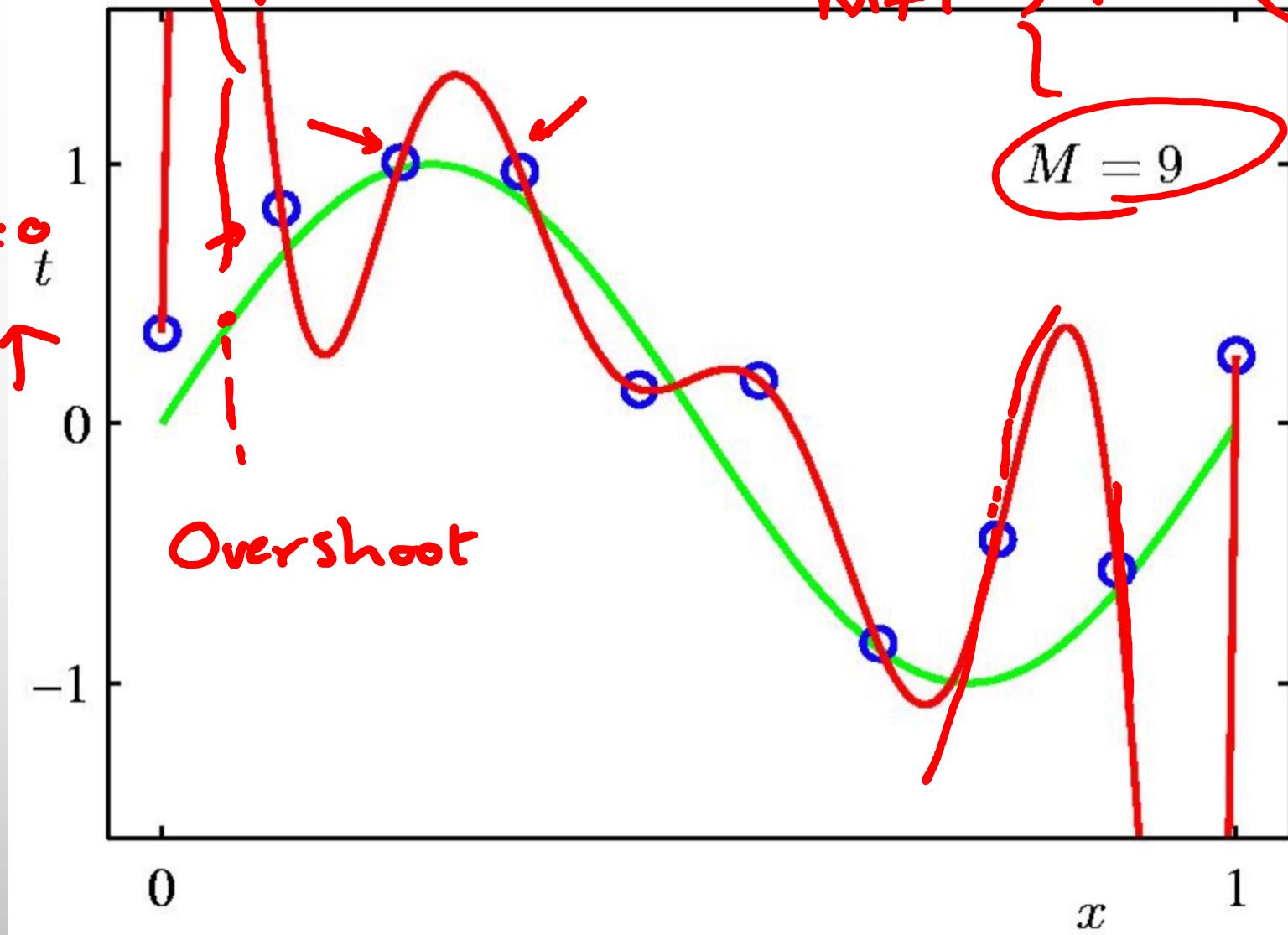


9th Order Polynomial

Overfit

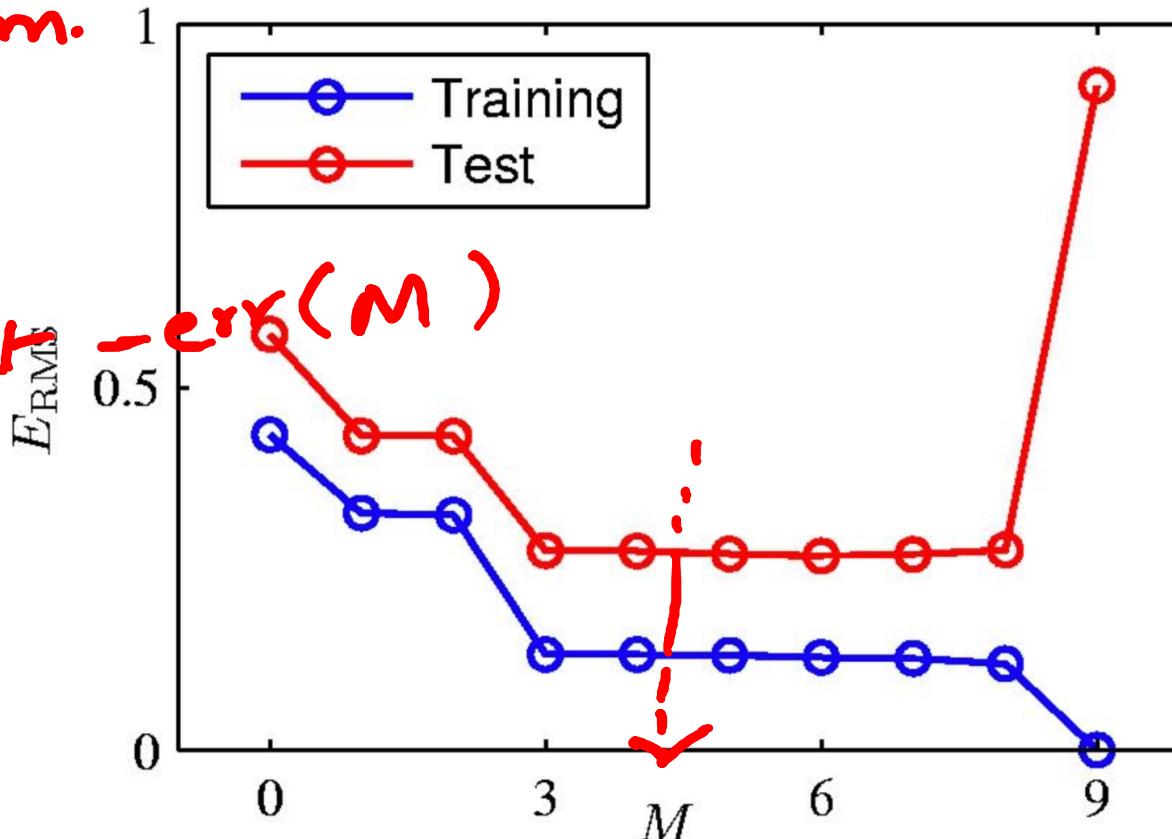
$$\underline{t}^M \underline{w} = \underline{0}$$

$$\text{1st err } \mathbb{T}$$



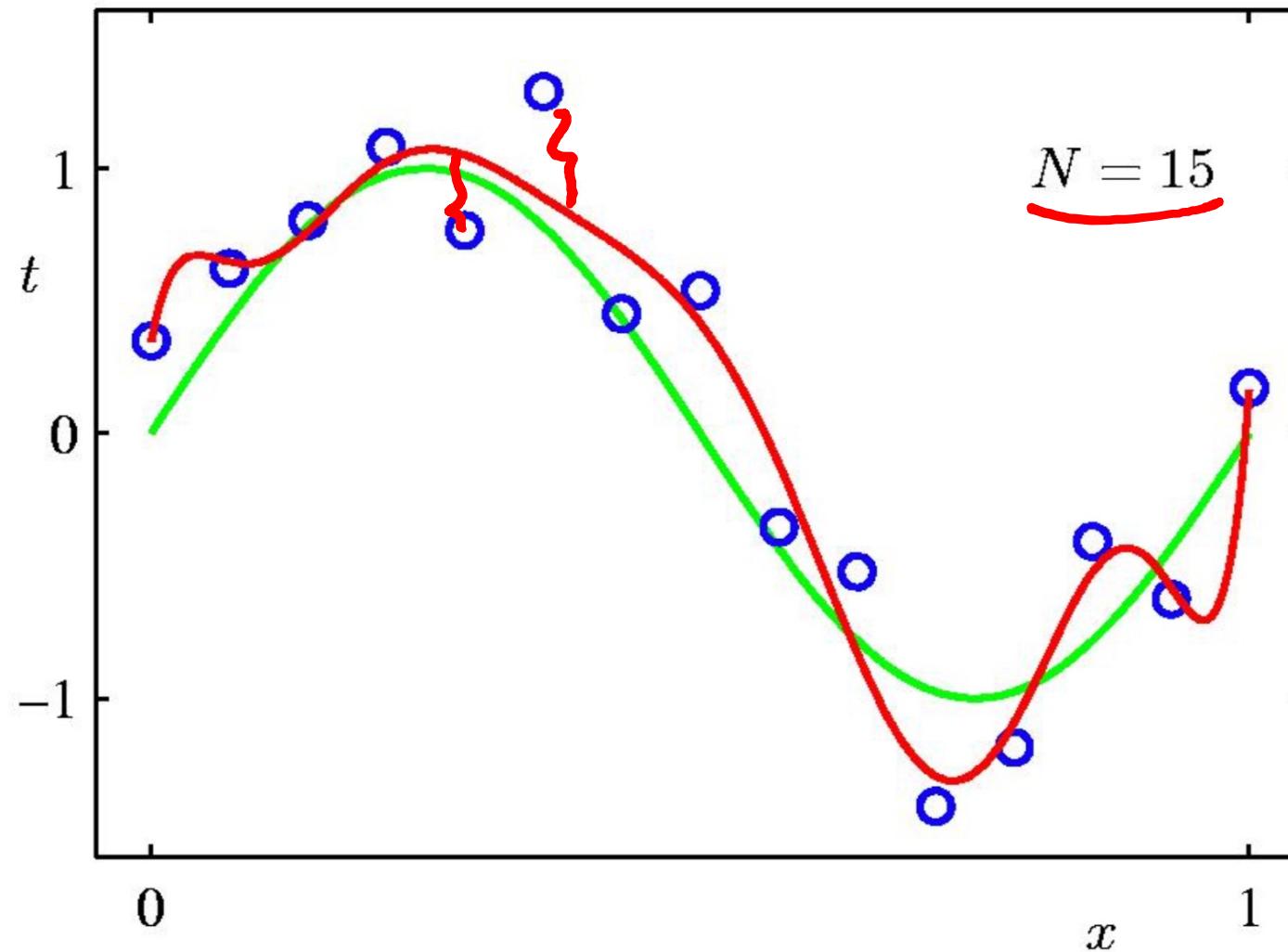
Over-fitting

Hyper param.
Tuning
min M Held-out E_{RMS} $\text{err}(M)$

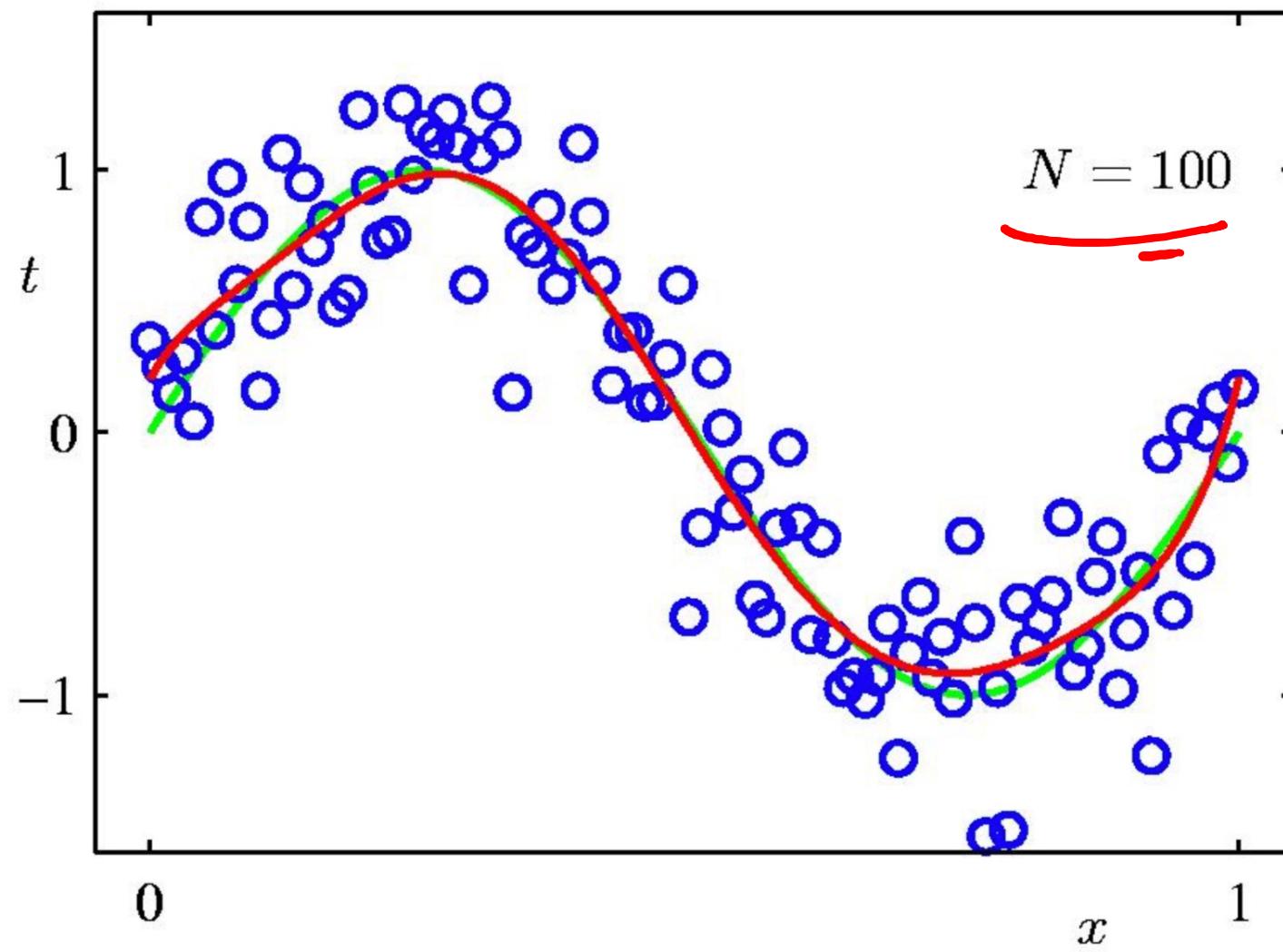


Root-Mean-Square (RMS) Error: $E_{\text{RMS}} = \sqrt{2E(\mathbf{w}^*)/N}$

Data Set Size: 9th Order Polynomial



Data Set Size: 9th Order Polynomial



training samples
 $N = 10$

Polynomial Coefficients

Overfit

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

$$f(\underline{x}) = \sum w_j x^j$$

$$\text{MSE} + \sum_j |w_j|$$

$$\text{MSE} + \lambda \sum w_j^2$$

Regularization

$$\underline{x} \rightarrow \begin{bmatrix} 1 & x & x^2 & \dots & x^m \end{bmatrix} \quad f = w^T \underline{\phi}$$

$$\| \nabla f \|^2 = \| w \|^2$$

Convex $(1, 5, 0, 0, 0, 0)$

Regularization

$$\widetilde{\|w\|}_1 = \|w\|_0$$

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

$$f = \underline{w^T x}$$

- Penalize large coefficient values

$$\frac{\lambda}{2} \widetilde{\|w\|}^2$$

regularization

ℓ_1 - regularization

Regularization:

$$\ln \lambda = -18$$

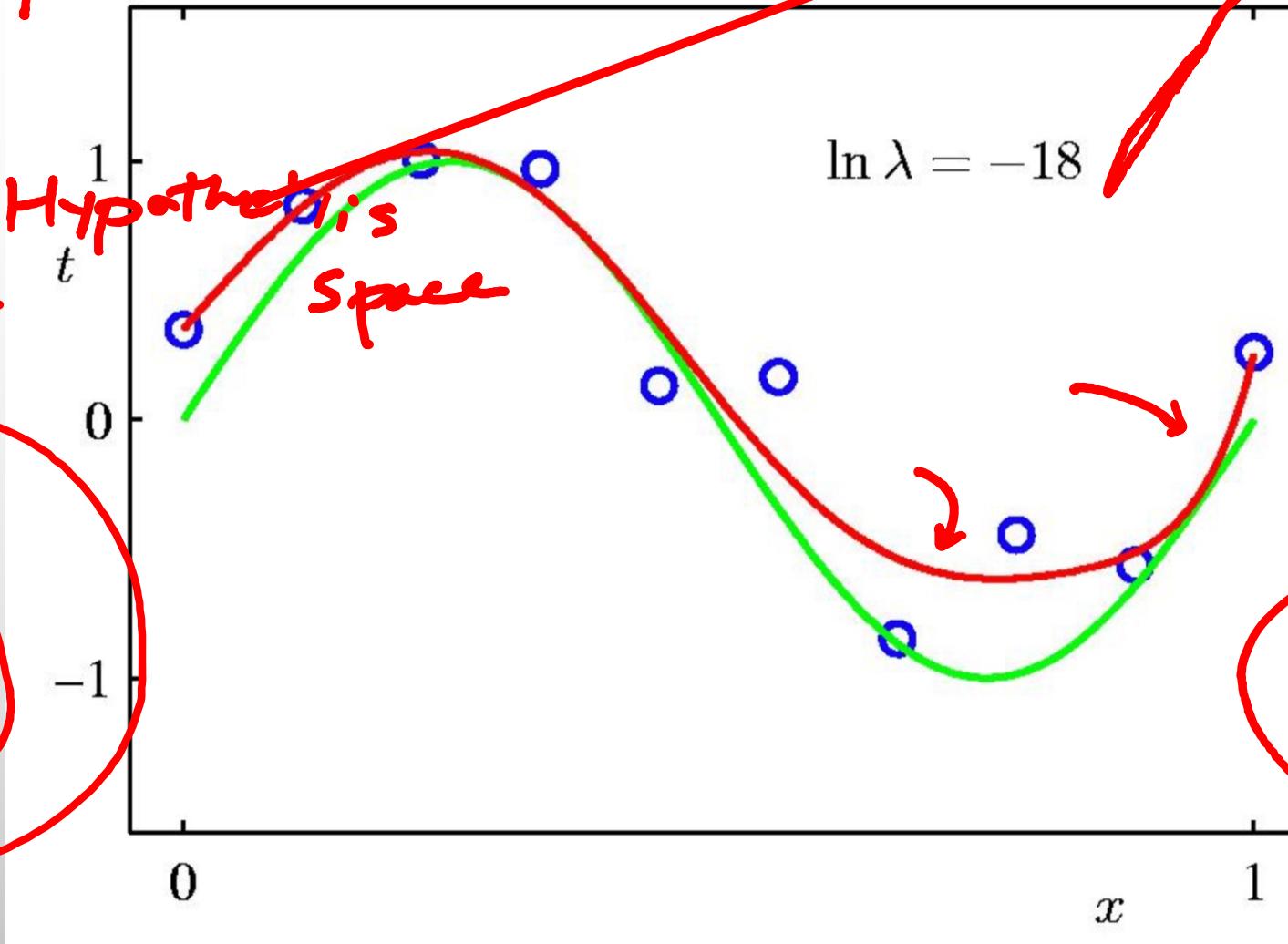
$M = 9$

$N = 10$

Hypothesis Space

$$\ln \lambda = -18$$

$M = 3$

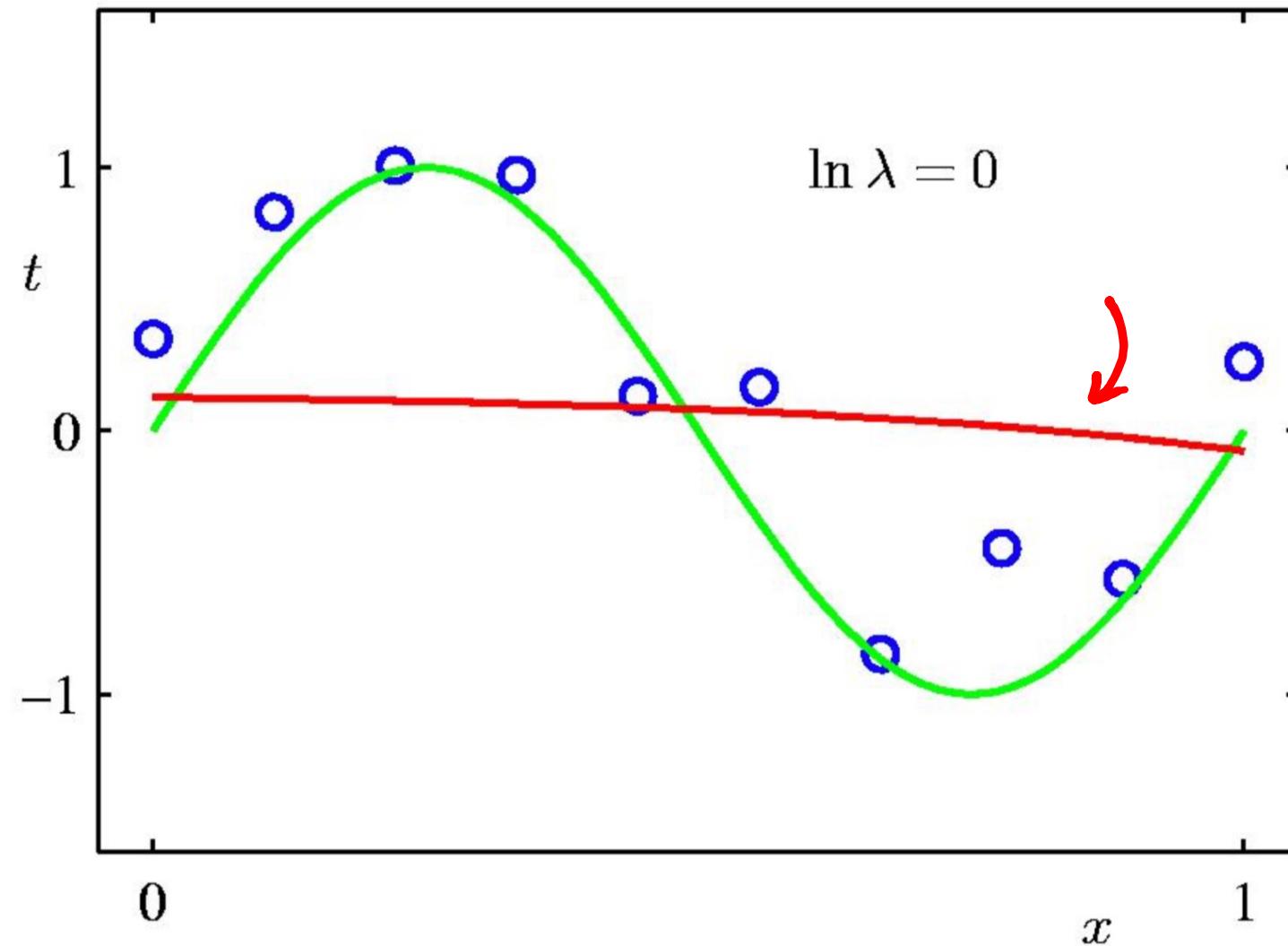


Regularization:

$$\ln \lambda = 0$$

$$\lambda = \underline{\underline{1}}$$

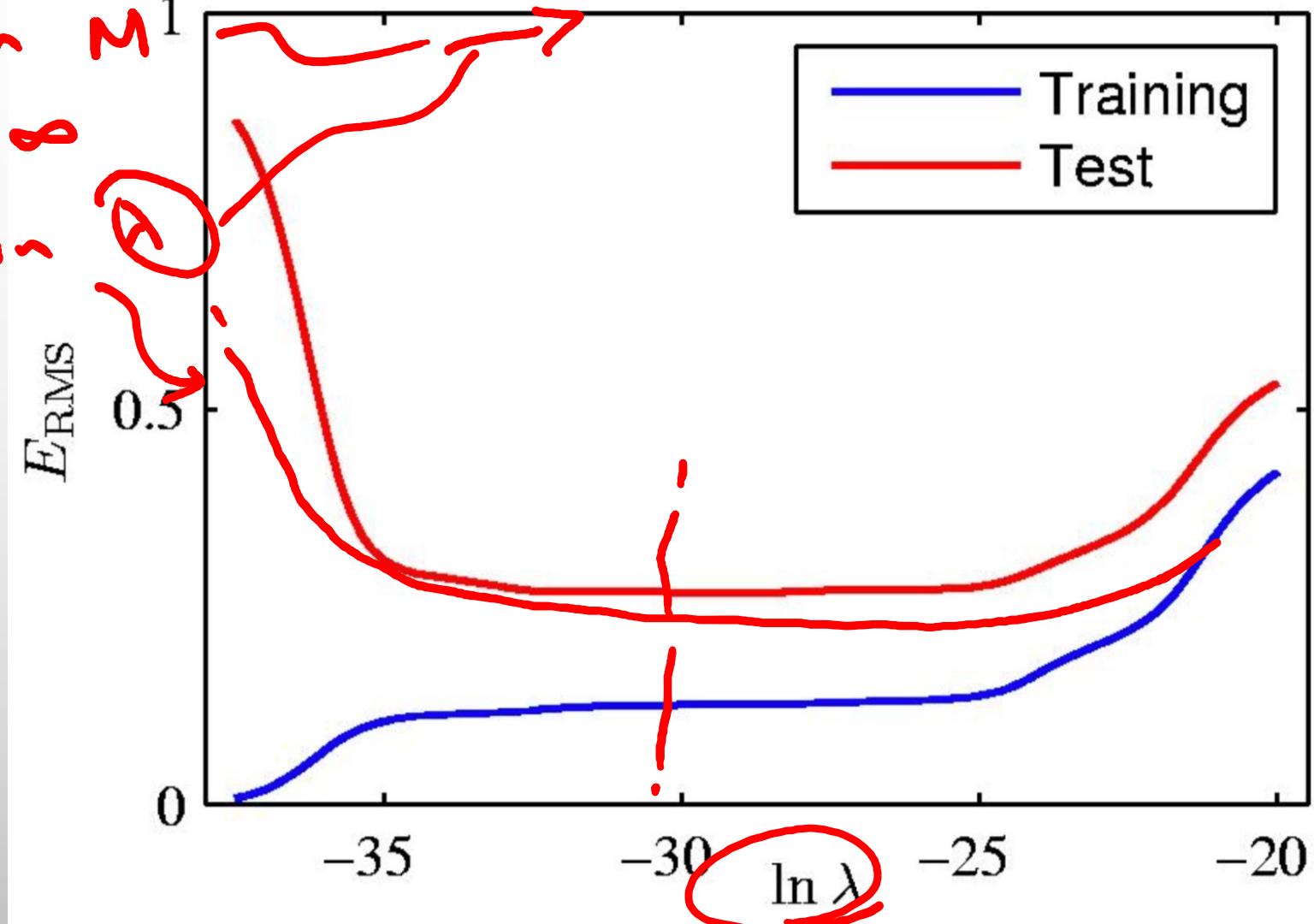
Underfit



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

which one is better ?

- HT on M^1
- $N \rightarrow \infty$
- HT on ∞



Polynomial Coefficients

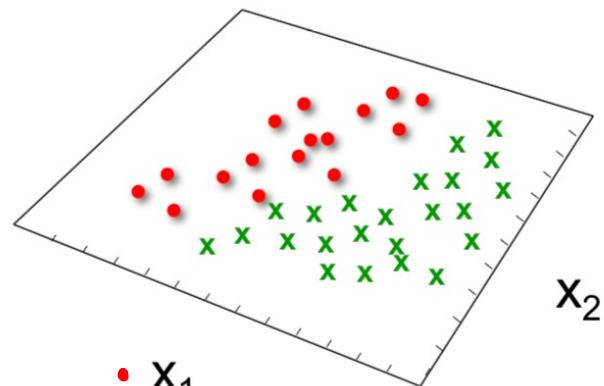
	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
w_0^*	0.35	0.35	0.13
w_1^*	232.37	4.74	-0.05
w_2^*	-5321.83	-0.77	-0.06
w_3^*	48568.31	-31.97	-0.05
w_4^*	-231639.30	-3.89	-0.03
w_5^*	640042.26	55.28	-0.02
w_6^*	-1061800.52	41.32	-0.01
w_7^*	1042400.18	-45.95	-0.00
w_8^*	-557682.99	-91.53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

Confidence Value

Logistic Regression $\xrightarrow{\text{tool}}$ Classification

$$\frac{P(Y = 1|X = \langle X_1, \dots, X_n \rangle)}{P(Y = 0|X = \langle X_1, \dots, X_n \rangle)} = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

implies



w is a
signature of

class 0

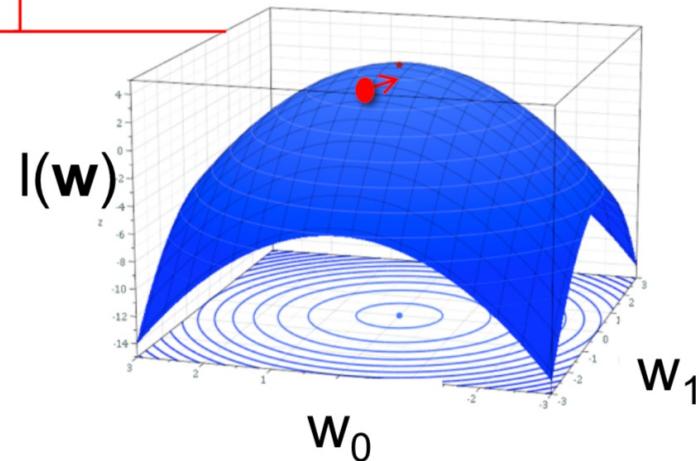
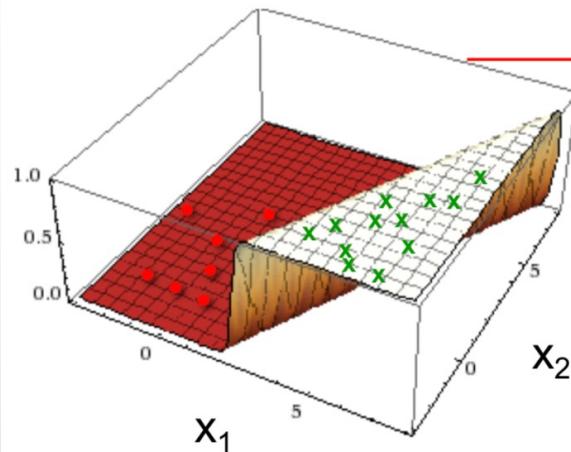
linear
classification
rule!

$$\ln \frac{P(Y = 0|X)}{P(Y = 1|X)} = w_0 + \sum_i w_i X_i$$

Gradient Ascent

$$w_0 = 40 \quad w_1 = -10 \quad w_2 = 5$$

Maximize $l(\mathbf{w}) = \ln P(D_Y | D_x, H_w)$



Update rule:
$$\Delta \mathbf{w} = \eta \nabla_{\mathbf{w}} l(\mathbf{w})$$

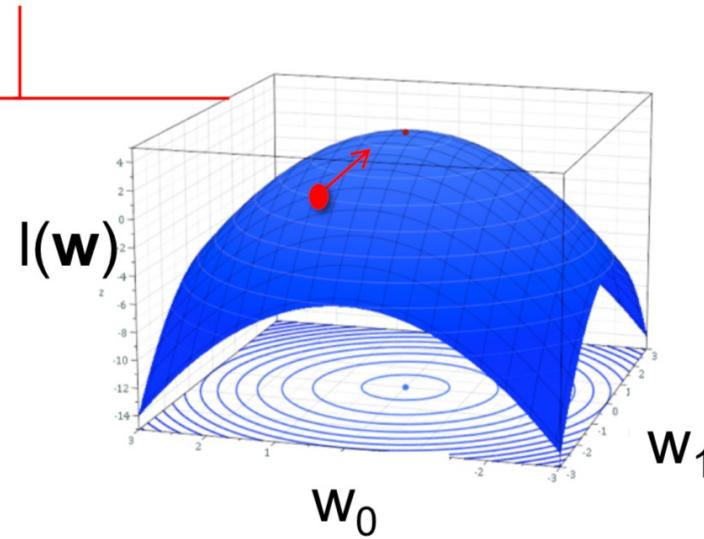
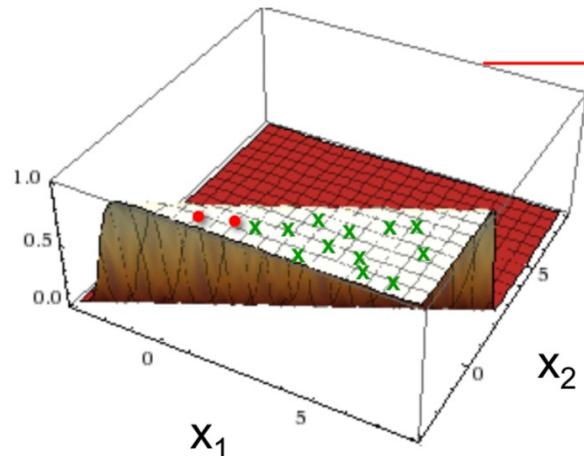
$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \frac{\partial l(\mathbf{w})}{\partial w_i}$$

Logistic w/ Initial Weights

$$w_0 = 20 \quad w_1 = -5 \quad w_2 = 10$$

$\text{Loss}(H_w) = \text{Error}(H_w, \text{data})$

Minimize Error \rightarrow Maximize $l(\mathbf{w}) = \ln P(D_Y | D_x, H_w)$



Update rule:

$$\Delta \mathbf{w} = \eta \nabla_{\mathbf{w}} l(\mathbf{w})$$

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \frac{\partial l(\mathbf{w})}{\partial w_i}$$

Step size

$$h = \sum \alpha_i x_i x_i^T \quad \alpha_i \geq 0$$

How to learn \underline{W} from the

training data $\mathcal{D} = \{(\underline{x}_1, y_1), \dots, (\underline{x}_N, y_N)\}$? \xrightarrow{iid}

$$\rightarrow P(\underline{Y} = \underline{y} | \underline{X}) = \frac{1}{1 + \exp\{\underline{y}^T \underline{W}^T \underline{X}\}}$$

Max. Likelihood

$$P(Y = -1 | \underline{x}) = \frac{\exp(\underline{W}^T \underline{x})}{1 + \exp(\underline{W}^T \underline{x})}$$

$$\begin{aligned} \max_{\underline{W}} P(Y^{(1)} = y_1, \dots, Y^{(N)} = y_N | \underline{x}_1, \dots, \underline{x}_N) &= \prod_{i=1}^N P(Y^{(i)} = y_i | \underline{x}_i) \\ &= \left(\prod_{i=1}^N \frac{1}{1 + \exp(y_i \underline{W}^T \underline{x}_i)} \right) \end{aligned}$$

GD on \underline{W}

$$\min_{\underline{W}} \left\{ \sum_{i=1}^N \log \left(1 + e^{\underline{y}_i \underline{W}^T \underline{x}_i} \right) \right\}$$

