

Artificial Intelligence

CE-417, Group 2

Computer Eng. Department

Sharif University of Technology

Falls 2020

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).

2

Regression

real-valued
(Output)

VS.

Classification
(K-classes)

all polynomials of order M

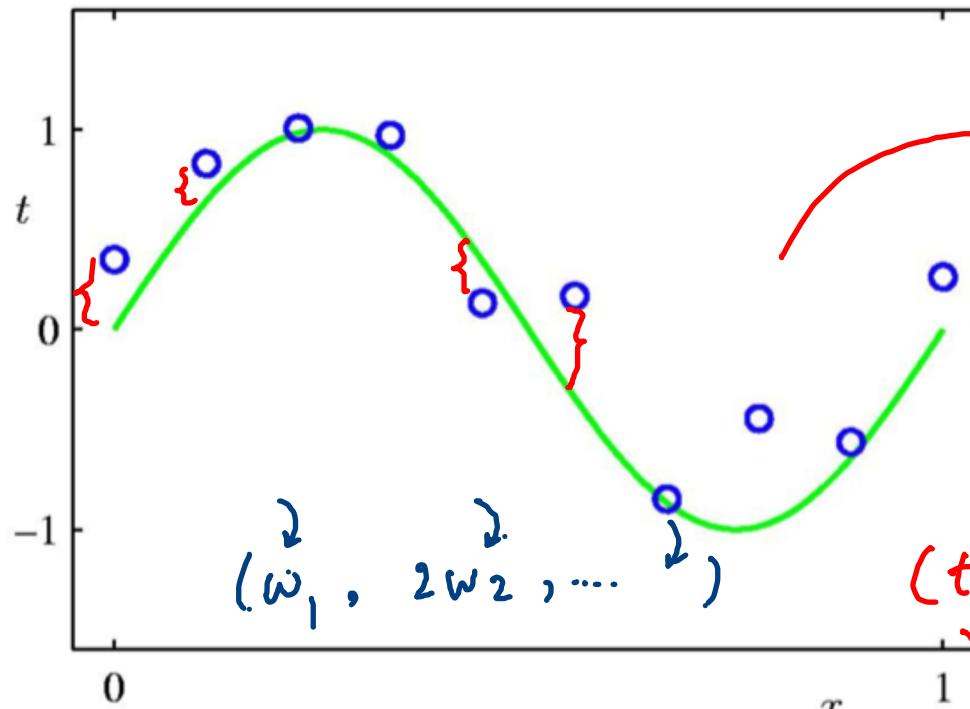
Polynomial Curve Fitting



$(w_0, w_1, \dots, w_{M+1})$

\mathbb{R}^{M+1}

Inductive Bias



training sample

$\text{Sin}(2\pi x)$

Search

Goal: The green
(target \leftarrow curve be
func.) estimated by

the blue points.

Hypothesis Space

Set of
all
candidate functions

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

No free Lunch \rightarrow determines the inductive bias

(trn. example)

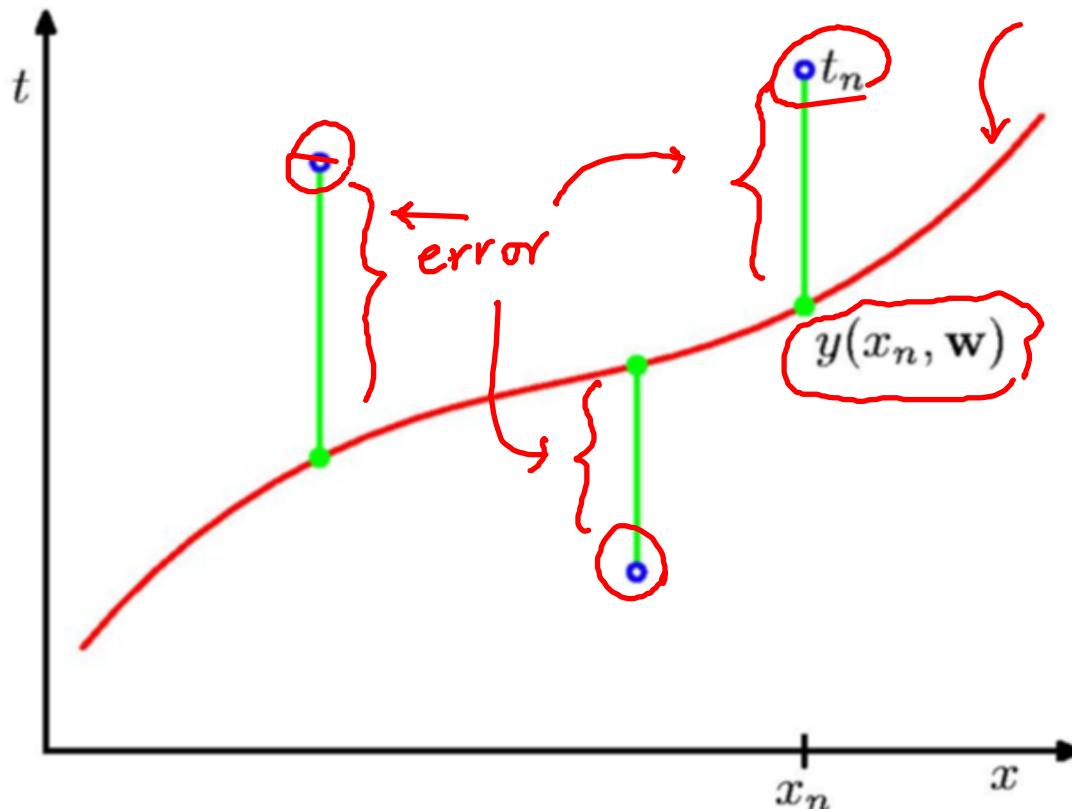
bias

How to
rank func.
in the fit?

Sum-of-Squares Error Function

$$\min_{w \in \mathbb{R}^M} E(w)$$

$$E(w) = \dots$$



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

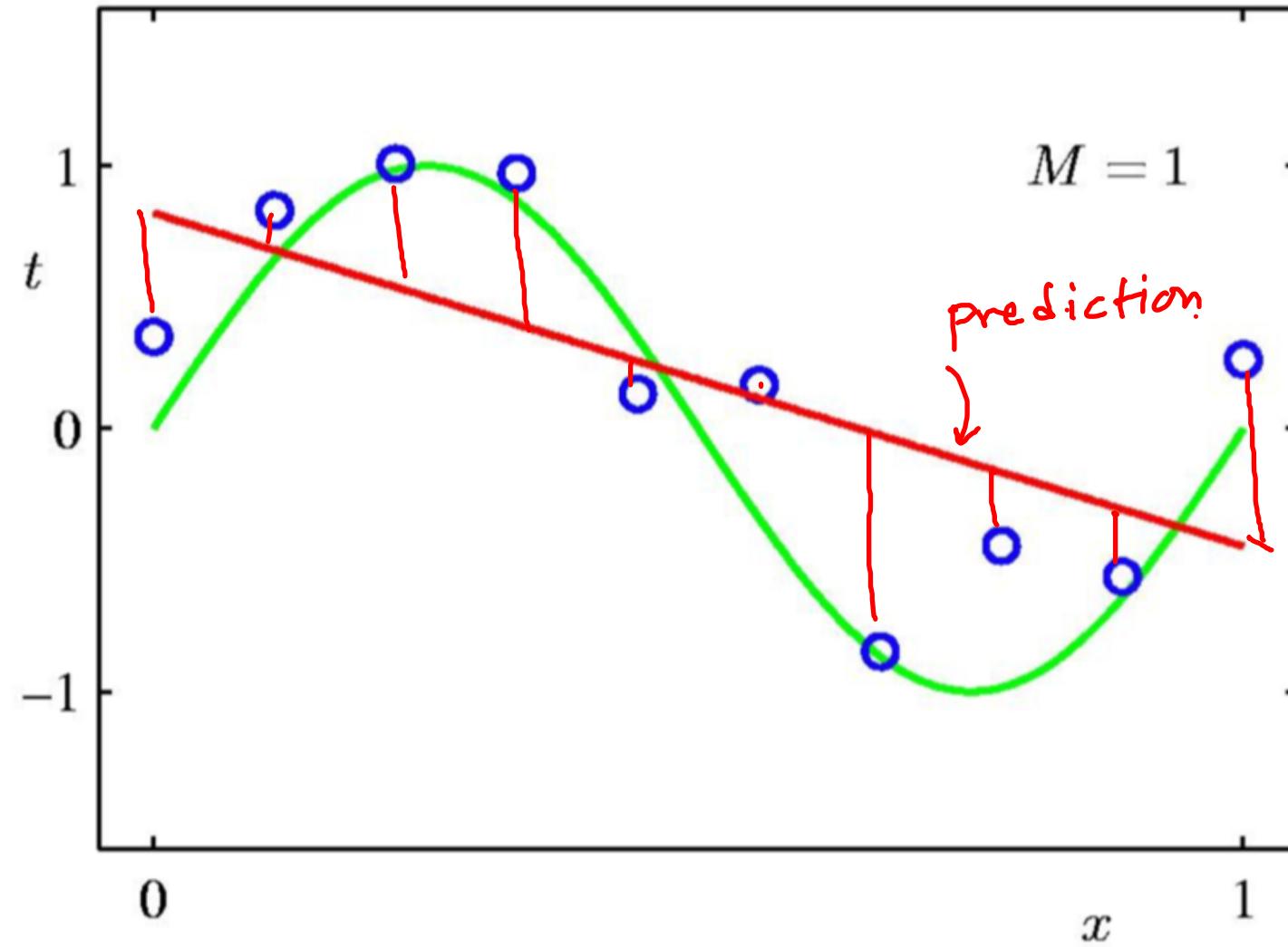
the func.
to be
ranked

n-th point

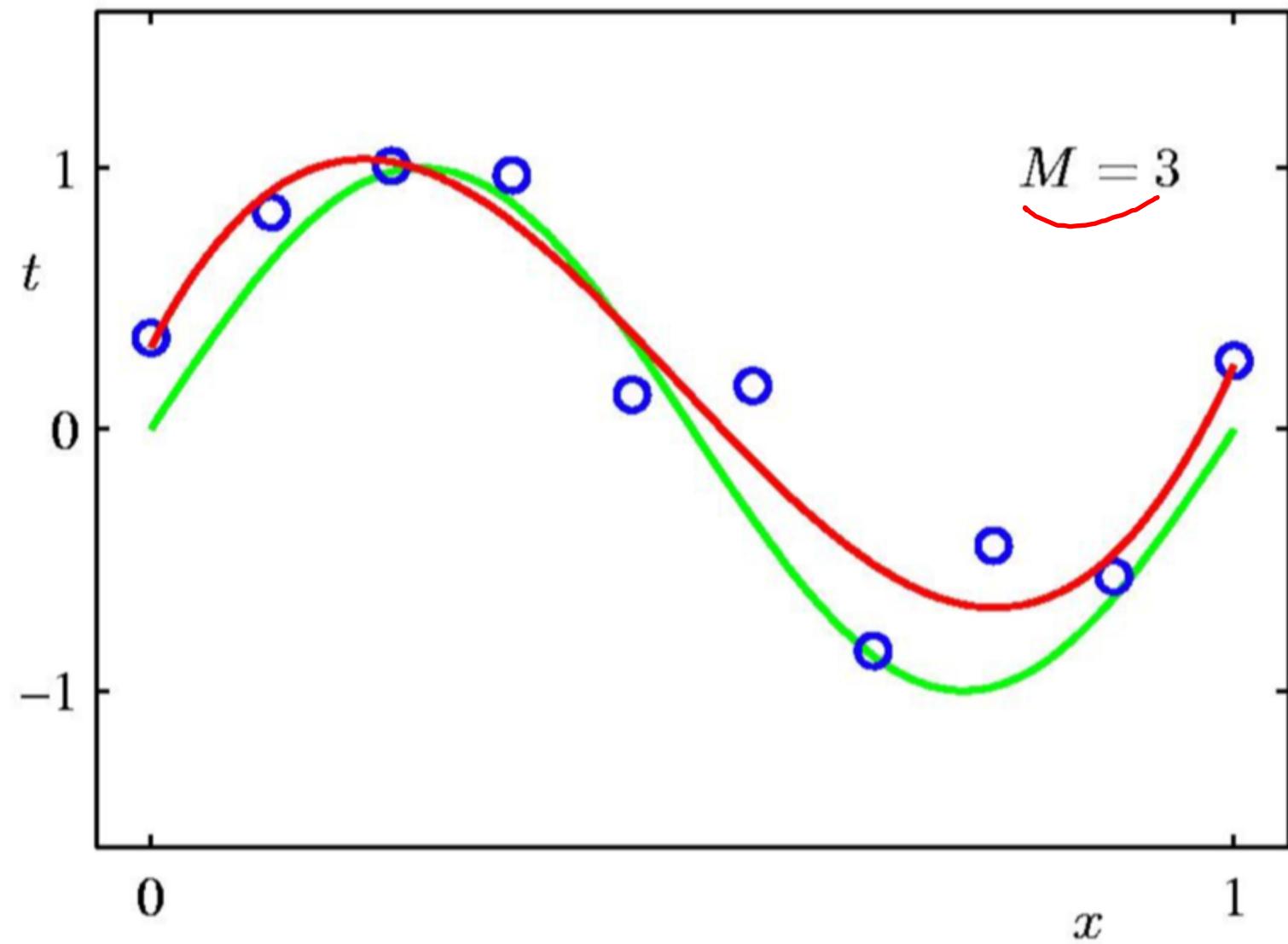
1st Order Polynomial

M = 1

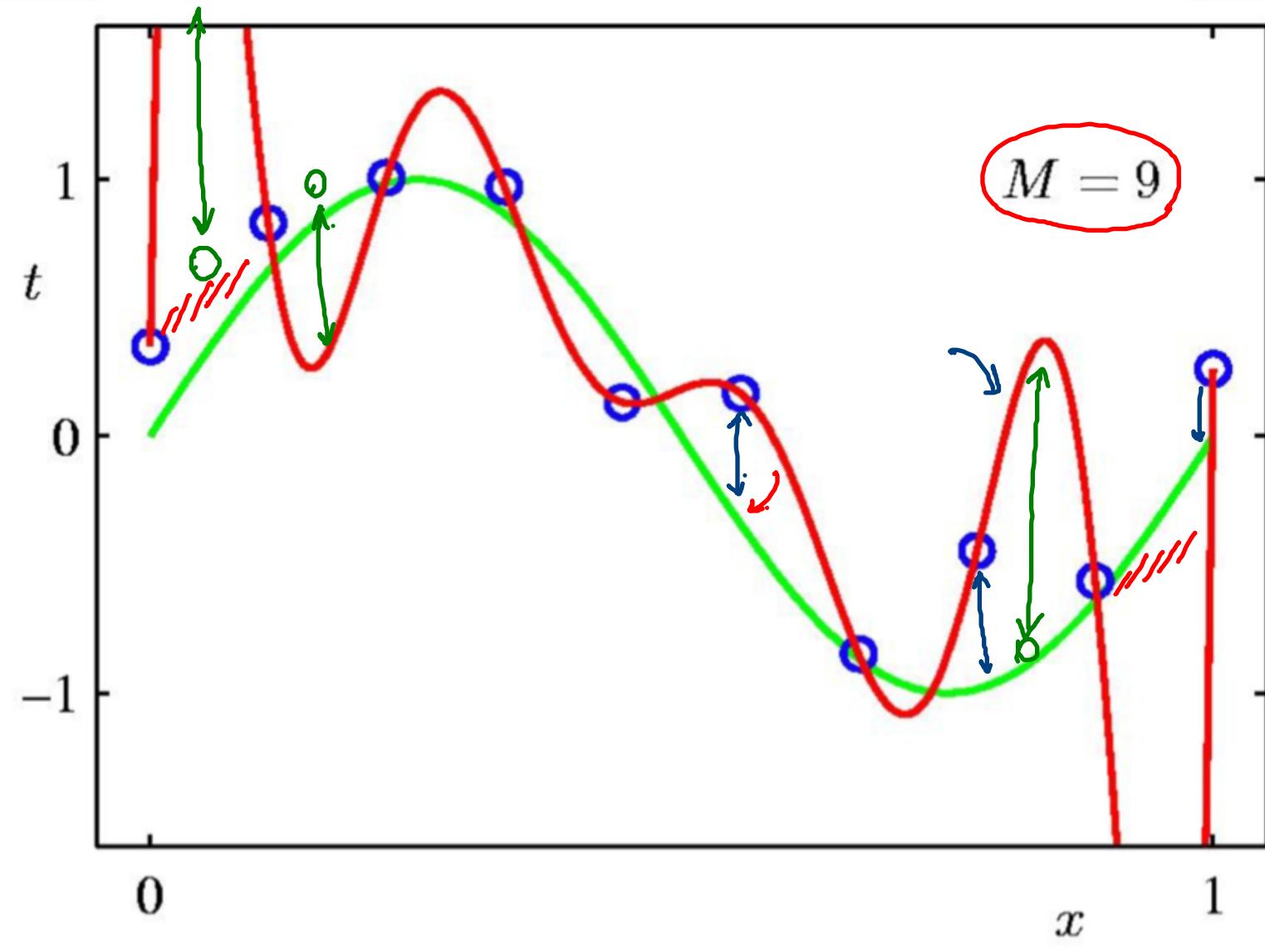
Gradient
Descent



3rd Order Polynomial



9th Order Polynomial



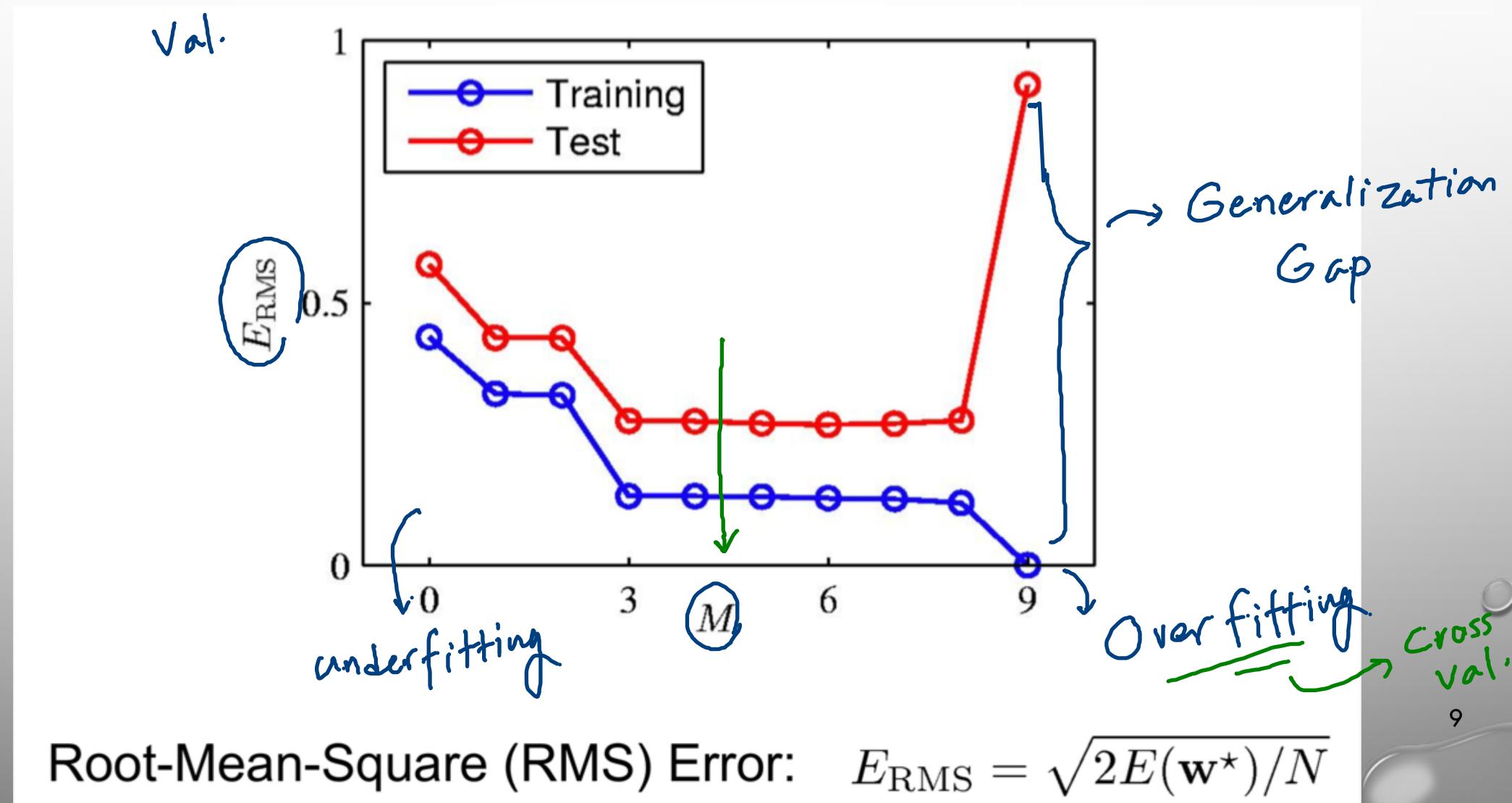
Training
error = 0

Test
error $\gg 0$

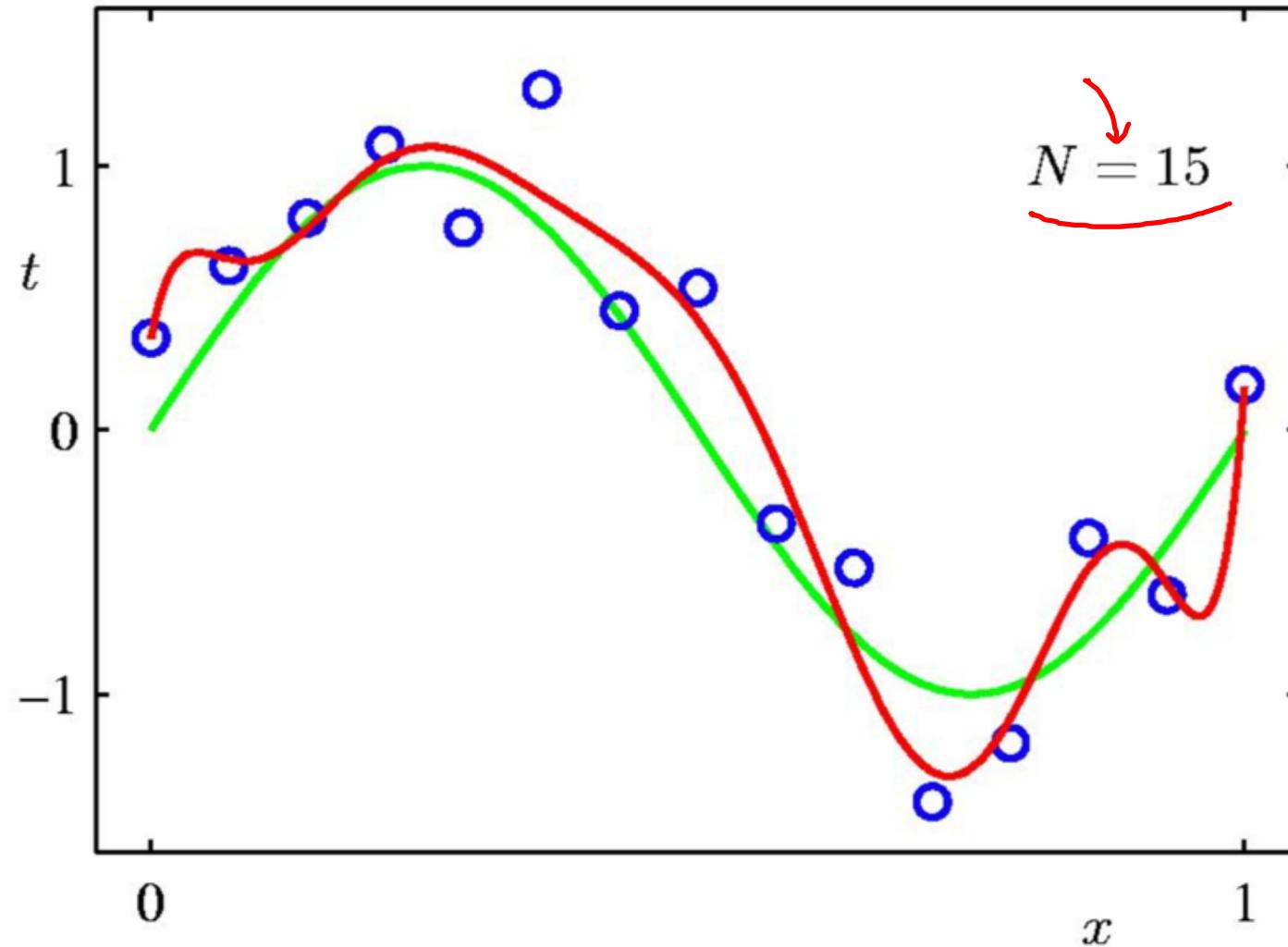
→ Over fitting

Over-fitting

P Chance Max



Data Set Size: 9th Order Polynomial



Data Set Size:

9th Order Polynomial

\Rightarrow Overfit is

not only

dependent

on the

Complexity

of our

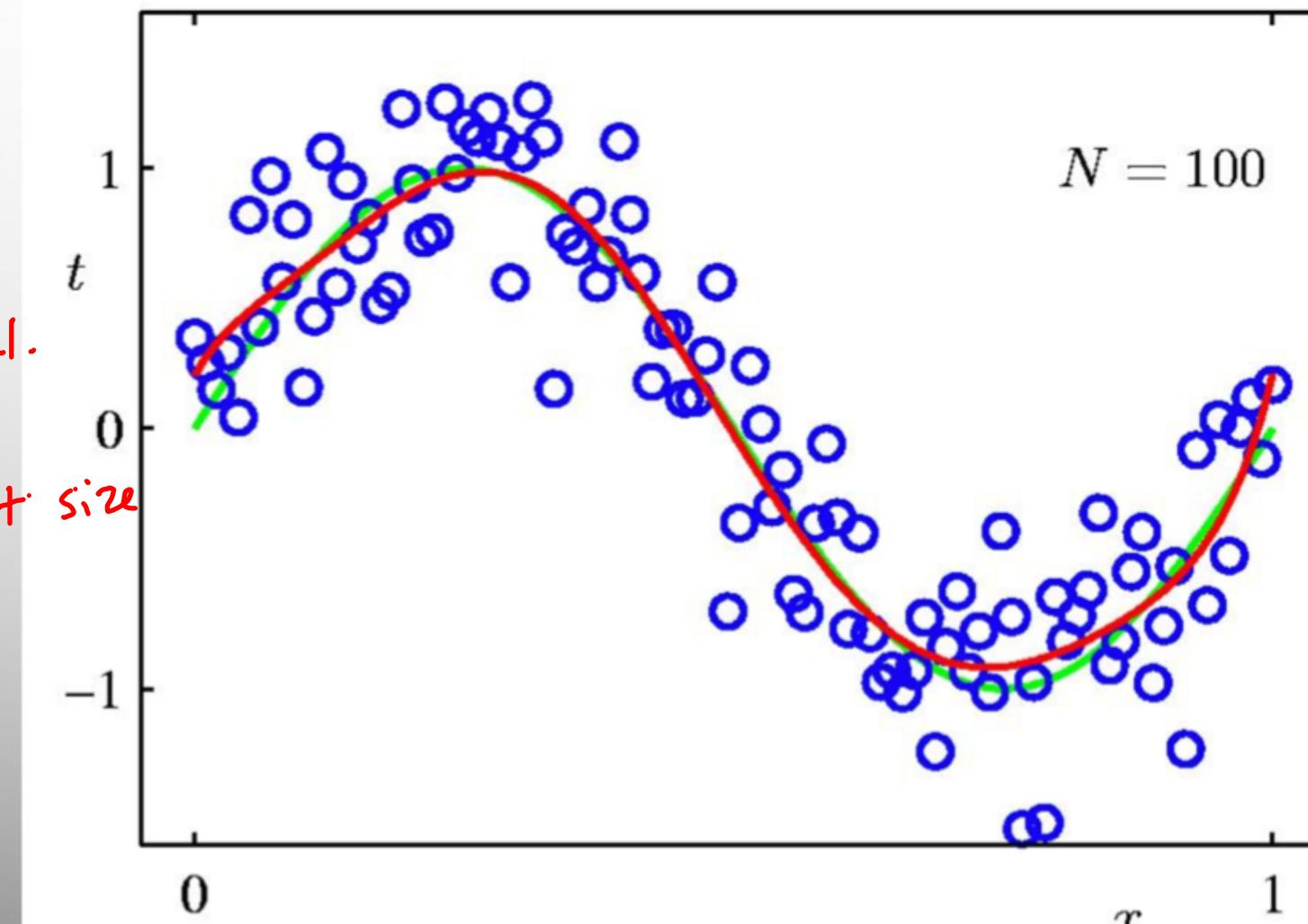
hypothesis

but also

depends on

trn. set size:

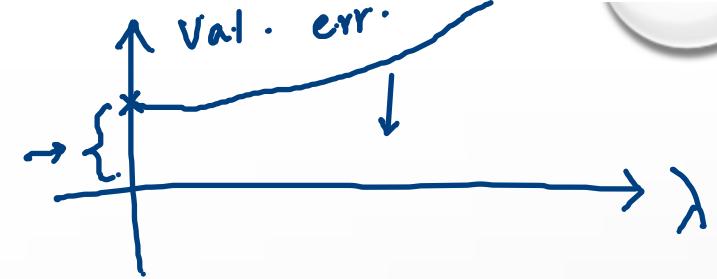
- ① Cross-val.
- ② Increase trn. set size



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

③ Regularization



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

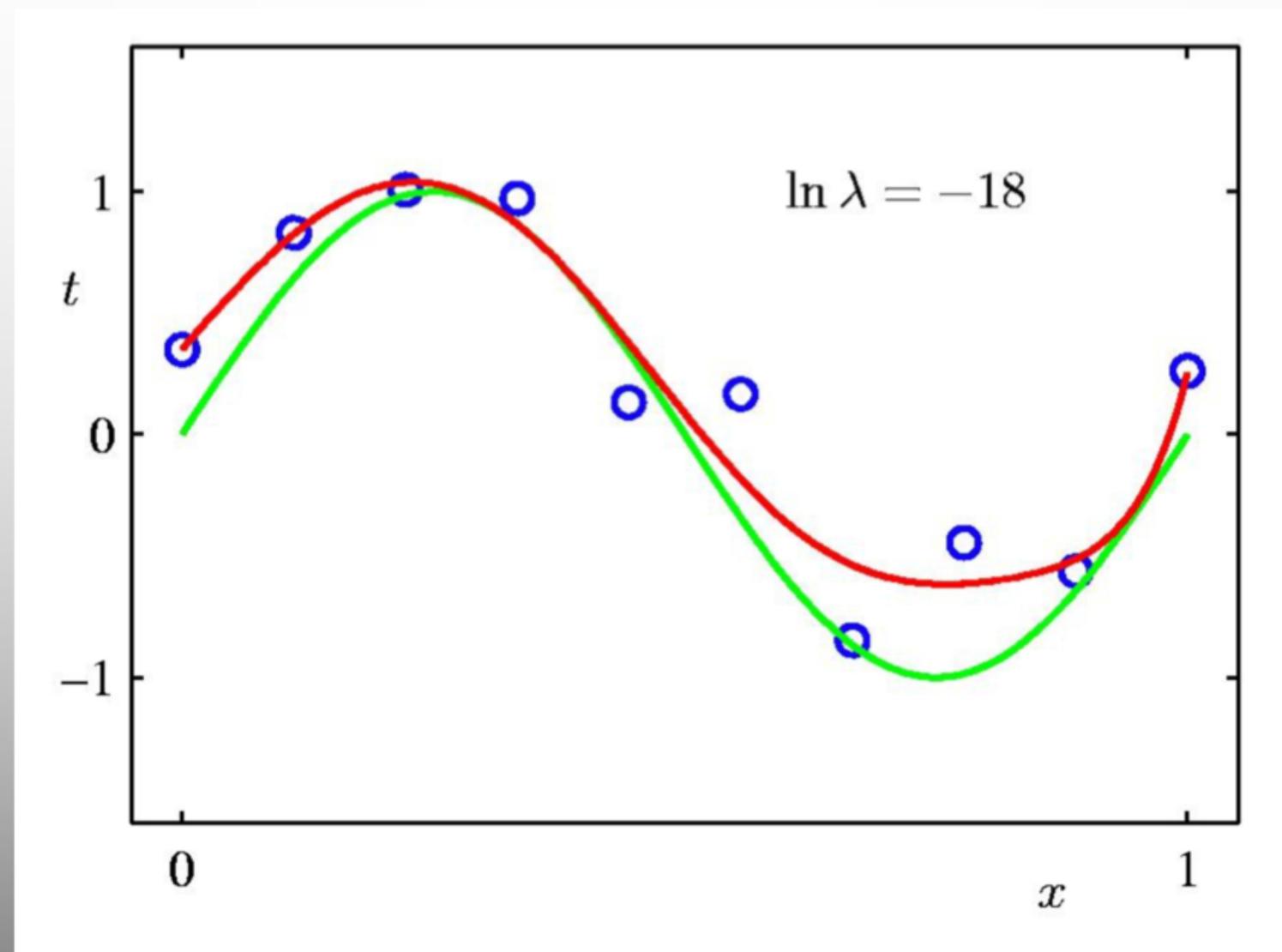
Empirical Err.

$\int_x \|\nabla_x y\|_2^2 dx$

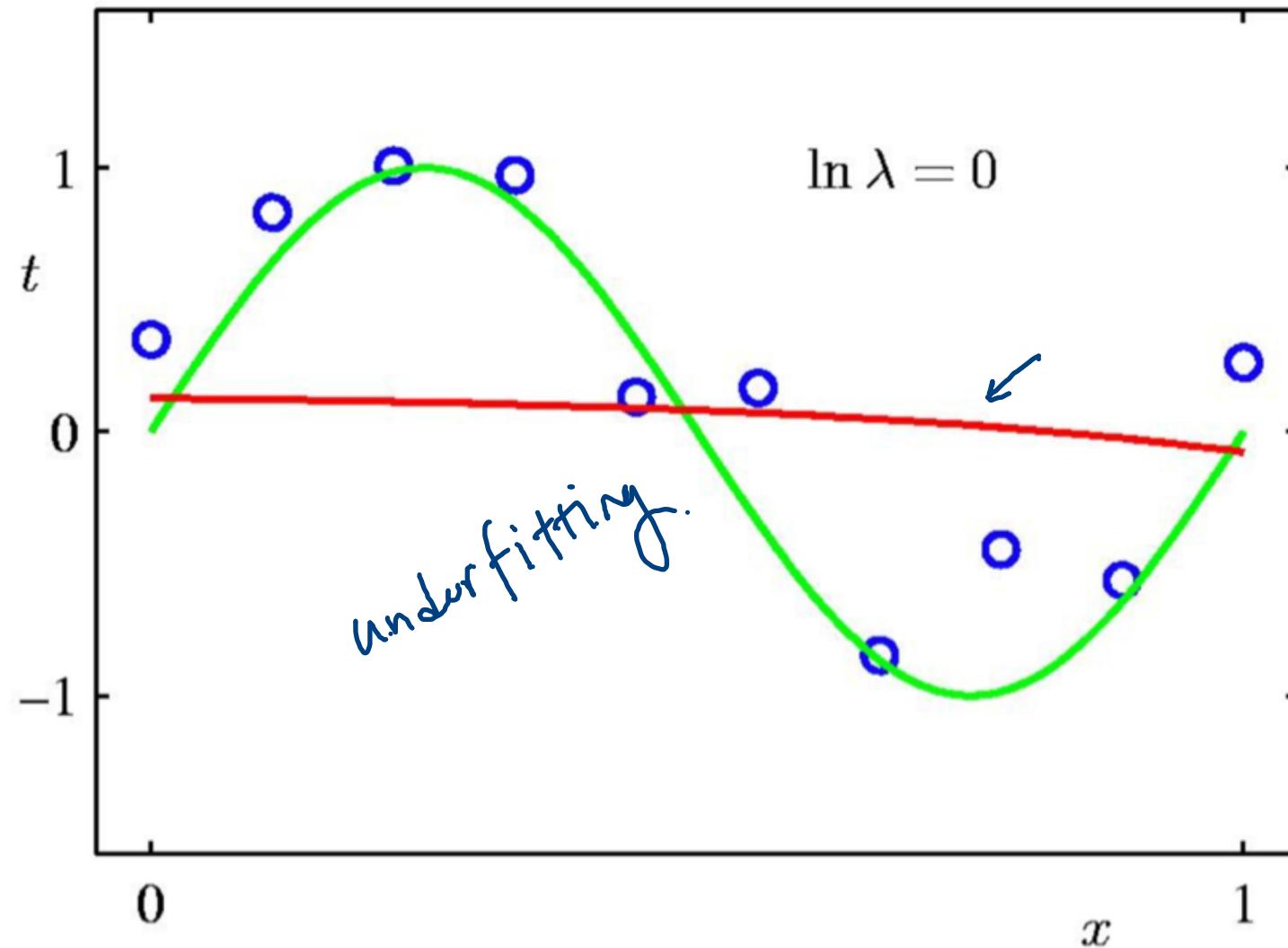
- Penalize large coefficient values

enforcing Smoothness \rightsquigarrow Derivative of any order
 ↳ the poly. Coeffs.

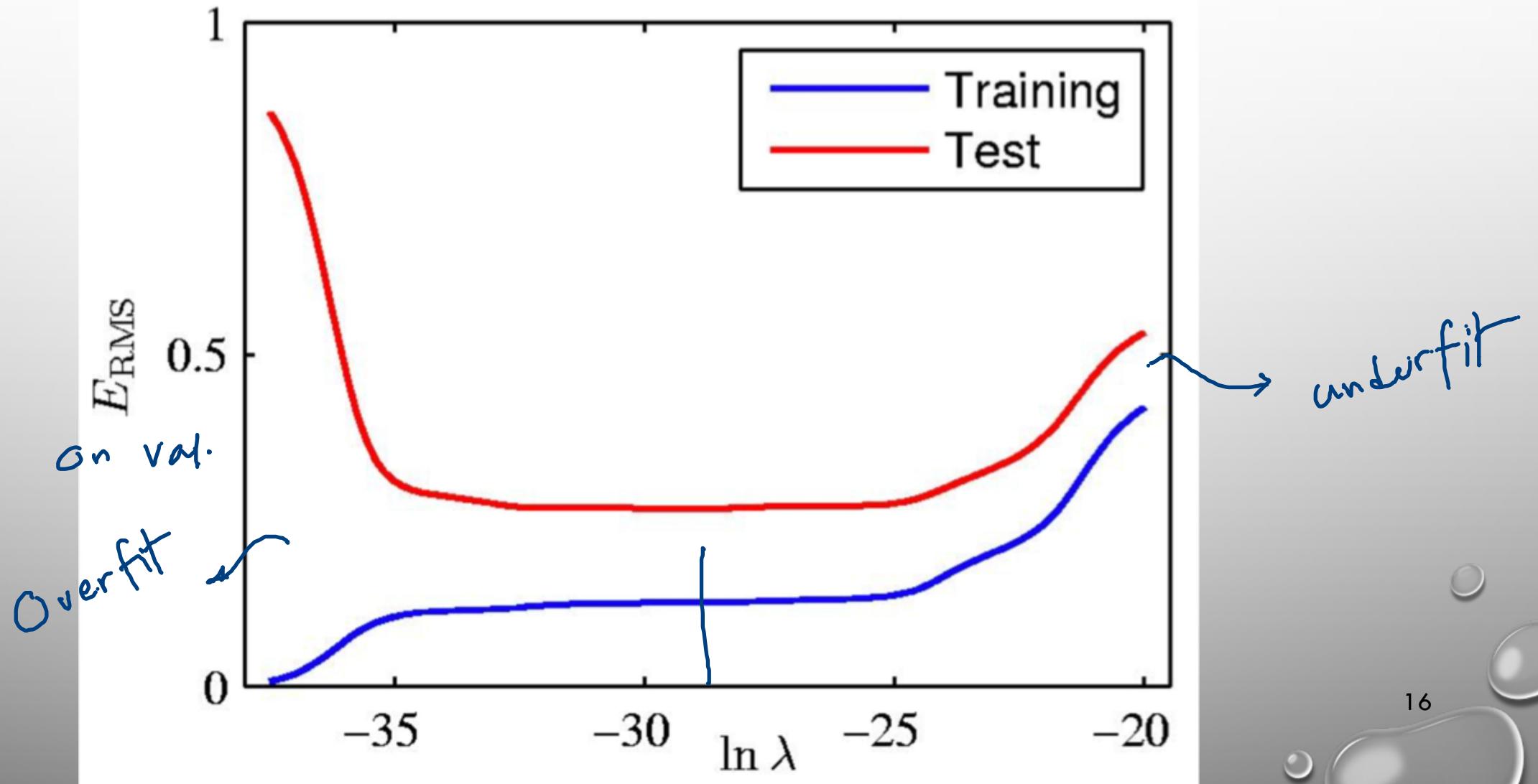
Regularization:

$$\ln \lambda = -18$$


Regularization:

$$\ln \lambda = 0$$


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



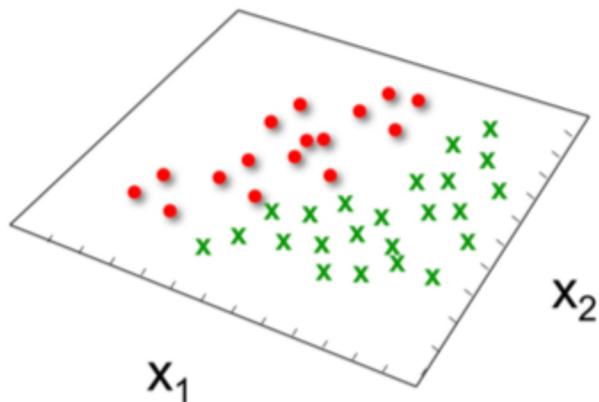
Polynomial Coefficients

	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
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

Logistic Regression

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

implies



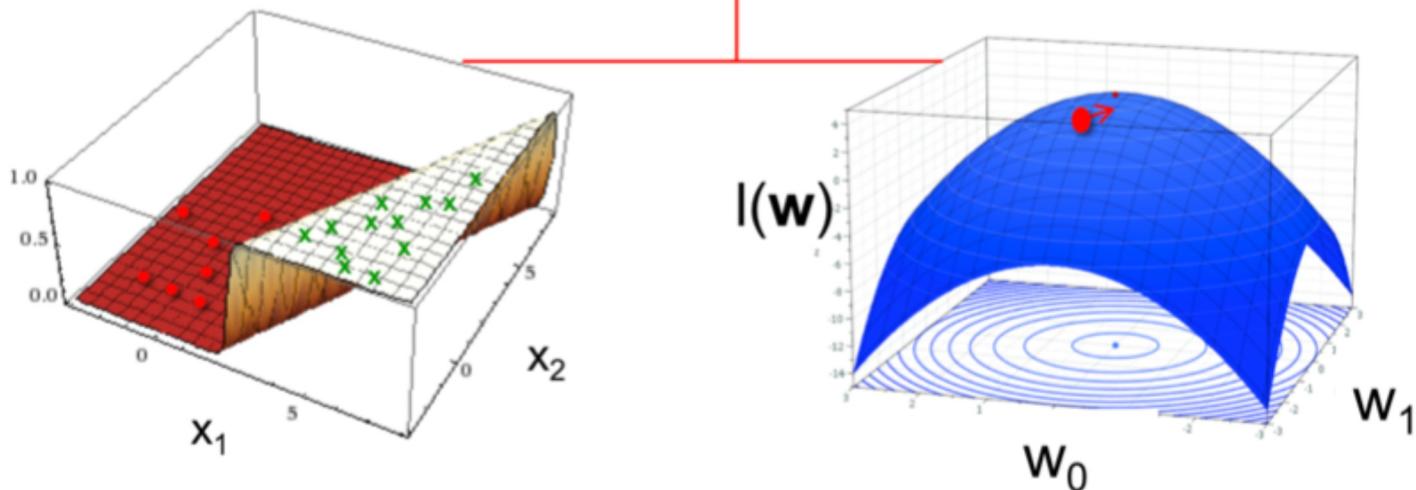
$$\ln \frac{P(Y = 0|X)}{P(Y = 1|X)} = w_0 + \sum_i w_i X_i$$

linear
classification
rule!

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

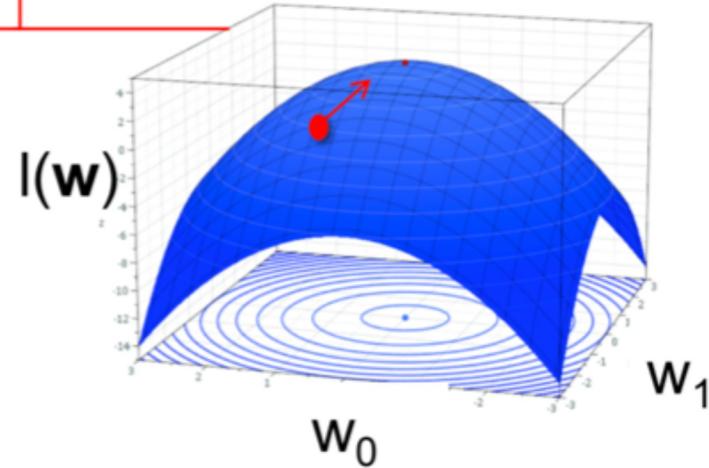
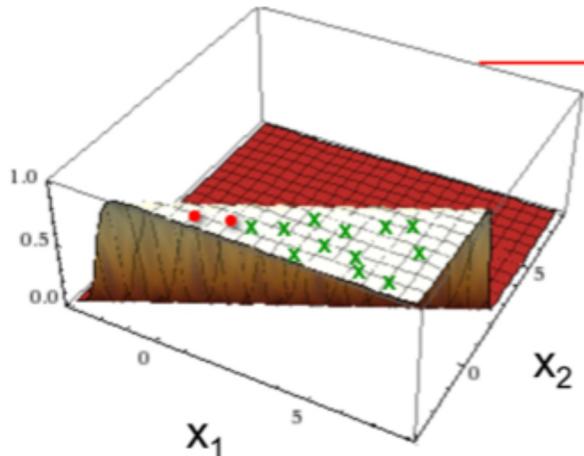
$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \frac{\partial l(\mathbf{w})}{\partial w_i}$$

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 $I(w) = \ln P(D_Y | D_x, H_w)$



Update rule:

$$\Delta w = \eta \nabla_w l(w)$$

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

Step size

x
 Weight Ch. Ins. ... \xrightarrow{y} Hearth Attack. \rightarrow Trn. set
 age Gender ... $(\underline{x}_1, 1), (\underline{x}_2, -1) \dots (\underline{x}_N, l)$

$$\underline{x} \in \mathbb{R}^d, \underline{w} \in \mathbb{R}^d$$

You can only predict its risk \rightsquigarrow real not the actual

Outcome
0/1

(prob. of the outcome)

$$P(Y=1 | \underline{x}, \underline{w}) = \frac{1}{1 + e^{-\langle \underline{w}, \underline{x} \rangle}} < 1$$

$$P(Y=-1 | \underline{x}, \underline{w}) = 1 - \frac{1}{1 + e^{-\langle \underline{w}, \underline{x} \rangle}}$$

$$= \frac{1 + e^{-\langle \underline{w}, \underline{x} \rangle} - 1}{1 + e^{-\langle \underline{w}, \underline{x} \rangle}}$$

$$P(Y=2 | \underline{x}, \underline{w}) = \frac{1}{(1 + e^{-\langle \underline{w}, \underline{x} \rangle})}$$

* What fitness func. to use in this scenario ?!

Max. Likelihood.

$$D = \{(\underline{x}_i, y_i) \mid 1 \leq i \leq N\} \xrightarrow{\text{iid}}$$

$$P(Y=y \mid \underline{x}, \underline{w})$$

→ which \underline{w} would make the prob. model
to assign the highest prob. to D ?

$$\max_{\underline{w}} P(D \mid \underline{w})$$

$$P(D \mid \underline{w}) = \prod_{i=1}^N P((\underline{x}_i, y_i) \mid \underline{w})$$

$$P((\underline{x}_i, y_i) \mid \underline{w}) = P(Y=y_i \mid \underline{x}_i, \underline{w})$$

$$P(\underline{x}_i \mid \underline{w})$$

$$\equiv \min_{\underline{w}} - \sum_{i=1}^N \log P(Y=y_i \mid \underline{x}_i, \underline{w})$$

$$\frac{1}{1 + e^{y_i \langle \underline{w}, \underline{x}_i \rangle}}$$

$$\equiv \min_{\underline{w}} \sum_{i=1}^N \text{log} \left(1 + e^{\frac{y_i \langle \underline{w}, \underline{x}_i \rangle}{A}} \right) \quad \begin{matrix} \rightsquigarrow \text{GD} \\ \rightsquigarrow \text{Convex.} \end{matrix}$$

L

$$\nabla_{\underline{w}} L = \sum_{i=1}^N \frac{y_i \underline{x}_i A}{1 + A} = \sum_{i=1}^N y_i \underline{x}_i \left(\frac{1}{1 + A^{-1}} \right)$$

$P(Y \neq y_i | \underline{x}_i, \underline{w})$

$$\underline{w}^{(t+1)} \leftarrow \underline{w}^{(t)} - \eta \sum_{i=1}^N y_i \underline{x}_i P(Y \neq y_i | \underline{x}_i, \underline{w})$$

if 1, the model did poorly on \underline{x}_i , so rotate \underline{w} in the \underline{x}_i dir. $\Rightarrow \underline{x}_i$ will not contribute to $\Delta \underline{w}$

if 0, the model has done well for \underline{x}_i

$$\begin{aligned}
 \nabla_w^2 L &= \nabla_w \left(\cdot \sum_{i=1}^N \frac{(y_i - x_i)^T}{1 + A^{-1}} \right) \\
 &= \sum_{i=1}^N \frac{\cancel{1/A^2} \cancel{y_i^2} \cancel{x_i x_i^T}}{\underline{(1+A^{-1})^2}} = \sum_{i=1}^N \underbrace{\left(\frac{C}{B}\right)}_{\alpha_i} \cancel{x_i x_i^T} \\
 &= \sum_{i=1}^N \alpha_i \cancel{x_i x_i^T} \quad \cancel{y_i} \\
 &\quad \text{rank } -1 \\
 &\quad \cancel{x_i} \\
 \end{aligned}$$

$\Rightarrow L$ is convex.

\Rightarrow ^{GD} on L converges to the global minimum
if η is suff. small.

