

# Machine Learning I: Supervised Methods

B. Keith Jenkins

## Announcements

---

- Slido event code: 5499248
  - slido.com
- Posted on D2L:
  - Class-participation grading criteria
- Announced on piazza:
  - Instructions to register regular recurring lecture-time conflicts
  - Homework file submission - please follow stated guidelines; otherwise you will lose points (HW4 and after)
- Homework 4 is due Friday
- No new homework will be posted this Friday
  - HW5 will be posted Friday 2/23
- Midterm exam is Wed., 3/6/2024
  - Vote for materials allowed: watch for piazza post with poll
  - 1 hour 45 min. exam
- No class on Monday 2/19 - Presidents' Day

## Reading

---

- Bishop 6.0-6.2, App. E
  - Lagrangian optimization; kernels

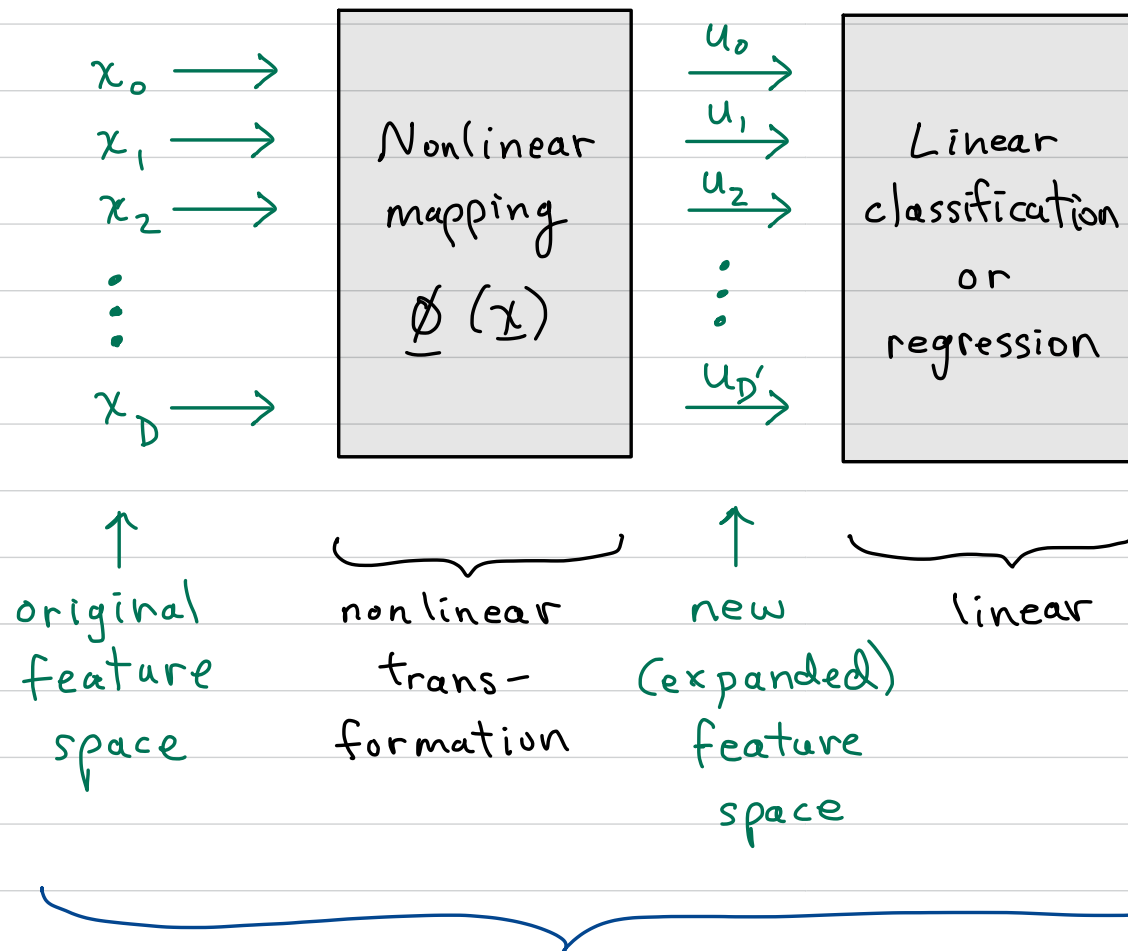
# Machine Learning I: Supervised Methods

B. Keith Jenkins

## Today's lecture

---

- Nonlinear classification and regression
  - Learning
- Complexity in supervised learning
  - Degrees of freedom
  - Constraints
  - Example



Classification:

$$g(\underline{u}) = \underline{w}_k'^T \underline{u}$$

$\rightarrow$  Linear in  $\underline{u}$

$\rightarrow$  Linear in  $\underline{w}_k'$

$\Phi$  machine (Nilsson) -  $\underline{\phi}(\underline{x})$  are polynomials in  $\underline{x}$  (originally)

Also called (for general nonlinear functions of  $\underline{x}$ ):

$\underline{\phi}(\underline{x})$ : "nonlinear transformation"  
or "basis-set expansion"

$g(\underline{u})$  or  $\hat{f}(\underline{u})$ : "generalized linear (discriminant or regression) functions"

## Learning nonlinear (in $\underline{x}$ ) models

1. Use nonlinear transformation  $\phi(\underline{x})$  to map  $\underline{x}$  into  $\underline{u}$ -space. (pre-processing)
- \* 2. Use any linear learning algorithm to find  $\hat{\underline{w}}'_k$  or  $\hat{\underline{w}}'$ .
3. Optionally, map:
  - (i) decision boundaries (and regions)
  - or (ii) regression function  $\hat{f}(\underline{u})$
 back into original feature space ( $\underline{x}$ -space).
4. How to make predictions on unknowns  $\underline{x}$ , given  $\hat{\underline{w}}'$  or  $\hat{\underline{w}}'_k$ ?
  - a.  $\underline{u} = \phi(\underline{x})$
  - \* b. Calculate  $g_k(\underline{u}) = \hat{\underline{w}}'^T_k \underline{u}$  then apply decision rule (classification)
  - \* or calculate  $\hat{y} = \hat{f}(\underline{u}) = \hat{\underline{w}}'^T \underline{u}$  (regression)
  - \* Same as linear case, except variable change  $\underline{x} \rightarrow \underline{u}$ ,  $\hat{\underline{w}} \rightarrow \hat{\underline{w}}'$

We will give other examples of  $\phi$  mappings later [with kernels SVM's].

# Complexity of Learning by Machine

Nonlinear transformation of $x$ before learning	# scalar weight variables to learn
--	---

---

Linear

$D+1$

2<sup>nd</sup> order  
polynomial

$\frac{1}{2}(D^2 + 3D) + 1$   
 $= O(D^2)$

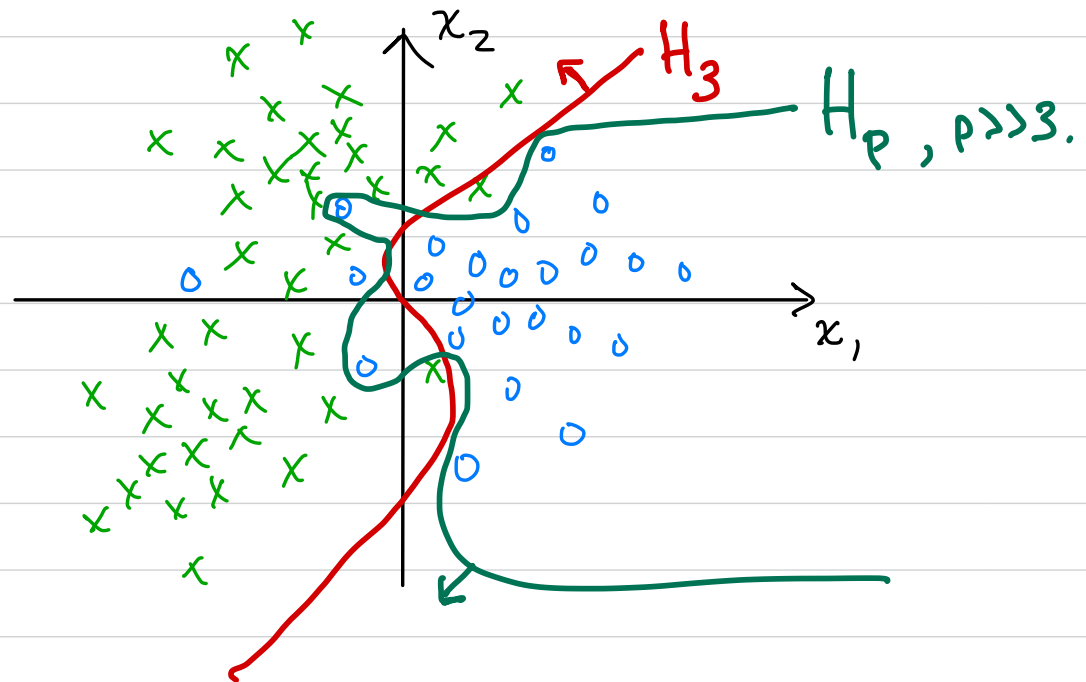
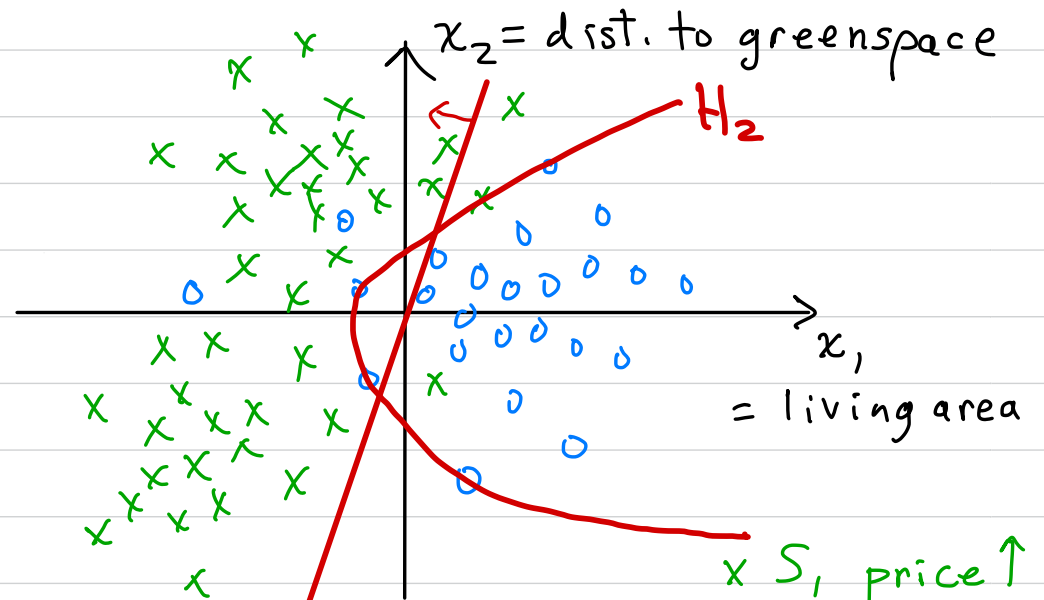
3<sup>rd</sup> order  
polynomial

$O(D^3)$

$p^{\text{th}}$  order  
polynomial

$O(D^p)$

⇒ Generalization depends on  
how flexible the decision  
boundary can be!



# Complexity of Learning - degrees of freedom (d.o.f.) and constraints ( $N_c$ )

Consider a system of linear equations and unknowns: (augmented notation)

$$\begin{array}{ccc} \underline{\underline{X}} & \underline{\underline{w}} & = \underline{\underline{b}} \\ \uparrow & \uparrow & \uparrow \\ [N \times (D+1)] & [D+1] & [N] \end{array} \quad (\text{orig. feat. space})$$

$$\text{or} \quad \begin{array}{ccc} \underline{\underline{\Phi}} & \underline{\underline{w'}} & = \underline{\underline{b'}} \\ \uparrow & \uparrow & \uparrow \\ [N \times (D'+1)] & [D'+1] & [N] \end{array} \quad (\text{expanded feature space})$$

#Eqs.	#unknowns (variables $w_j$ or $w'_j$ )	
$N$	$D+1$	(orig. feat. space)
$N$	$D'+1$	(expanded feat. space)

$\uparrow$   
 $\approx N_c$

$\uparrow$   
d.o.f.

We can define degrees of freedom (d.o.f.) in a learning problem as the number of variables that the learning algorithm can adjust independently.

Ex:  $g(\underline{x}) = \underline{w}^T \underline{x}$  (augm. notation) (2-class)

[  $\text{d.o.f.} = D+1$

Constraints constrain choices of these variables.

Ex:  $\underline{X} \underline{w} = \underline{b}$  is a set of  $N$  constraints on  $\underline{w}$ .

Let  $N_c = \# \text{ of constraints.}$

# Complexity: $N_c$ and d.o.f.

Consider:  $\underline{\underline{X}} \underline{\underline{w}} = \underline{\underline{b}}$  (original feat. spc.)

$\underline{\underline{\Phi}} \underline{\underline{w}}' = \underline{\underline{b}}'$  (expanded feat. spc.)

<u>Feature space</u>	<u><math>N_c \leftrightarrow</math> d.o.f.</u>	<u>system of eqns.</u>	<u>Comments</u>
orig.	$N = D+1$	exactly determined	} if <u><math>X</math></u> or <u><math>\Phi</math></u> is full rank
expanded	$N = D'+1$	" "	
(more generally, $N_c = \text{d.o.f.}$ )			
orig.	$N > D+1$	over determined	} if <u><math>X</math></u> or <u><math>\Phi</math></u> has full column rank.
expanded	$N > D'+1$	" "	
(more generally, $N_c > \text{d.o.f.}$ )			
orig.	$N < D+1$	under determined	}
expanded	$N < D'+1$	" "	
(more generally, $N_c < \text{d.o.f.}$ )			



Ex:Linear model for  
2-class problem $N=3$ 

(S)

$$g(\underline{x}) = w_0 + w_1 x_1 + w_2 x_2$$

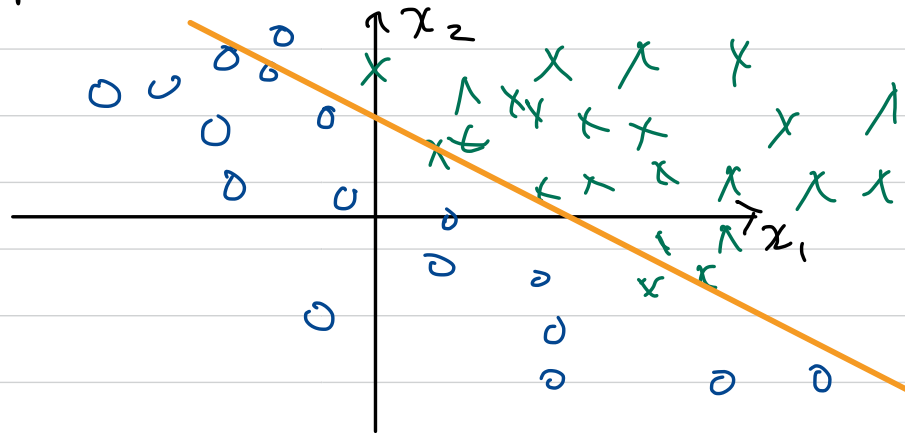
$$\underline{X} \underline{w} = \underline{b}$$

$$N_c = \# \text{ eqns.} = 3$$

$$\text{d.o.f.} = \# w_j = 3$$

 $\Rightarrow$  exactly determined if  $\underline{X}$  is of full rank. $\rightarrow$  In ML, we prefer to learn in overdetermined realm.

e.g.



$xS_1$  (house price  $\uparrow$ )  
 $\circ S_2$  (house price  $\downarrow$ )

wide variety of decision boundaries is possible.

General rule of thumb:

Typically we want  $N_c > (3-10) \times \text{d.o.f.}$

e.g.,  $N > (3-10) \times (D+1)$  for linear discr.fcn  
case, orig. feat. space,  
given above

(Empirical / numerical result.)

What if  $N_c < (3-10) \times \text{d.o.f.}$ ?

$\Rightarrow$  Risk of overfitting

Possible approaches:

1. Increase  $N_c$  by increasing  $N = \# \text{ data points}$

- Collect more data

- Synthetically generate more data

2. Increase  $N_c$  by adding other constraints

- Use prior knowledge in some way

- Restrict choices of optimal  $\underline{w}$

$\rightarrow$  3. Decrease d.o.f.

- Reduce  $D$  or  $D'$  (e.g., feature selection or transformation)

- Reduce amount of nonlinearity in the model



Today

## Example of d.o.f. and $N_c$

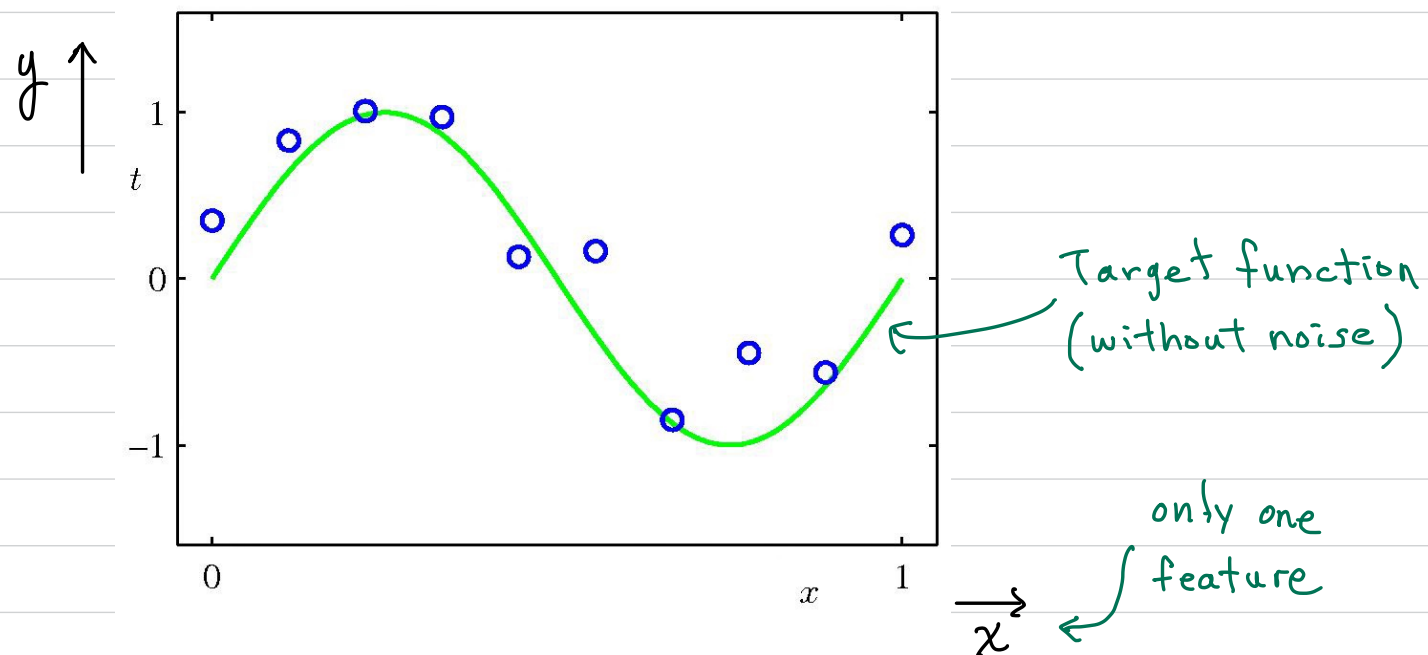
Consider a 1D regression problem [Bishop, Sec.1.1]

Let the known output values  $y_i$  in the dataset come from a target fcn:

$$y = \sin(2\pi x) + \text{noise}$$

$$\text{thus: } \mathcal{D} = \{(x_i, y_i), 0 \leq x_i \leq 1 \text{ (uniformly spaced on } [0, 1]), \\ y_i = \sin(2\pi x_i) + n_i, n_i \sim N(0, \sigma^2), i = 1, 2, \dots, N\}$$

Let  $N=10$ .



We will use a set of (mostly nonlinear) models - polynomials in  $x$ :

$$\hat{f}_d(x) = \sum_{i=0}^d w_i' x^i, \quad d \text{ is a parameter we specify.}$$

$$\text{Let } \underline{u} = \begin{bmatrix} x^0 \\ x^1 \\ x^2 \\ \vdots \\ x^d \end{bmatrix}, \quad \underline{w}' \text{ is defined above.}$$

⑤

Expanded f.s. dimensionality = ?

$d+1$

Original f.s. dimensionality = ?

1 feature  $x$

$\Rightarrow$  dim. 1 (non-augmented)

d.o.f. = ?

$d+1$  scalar  $w_j'$  values.

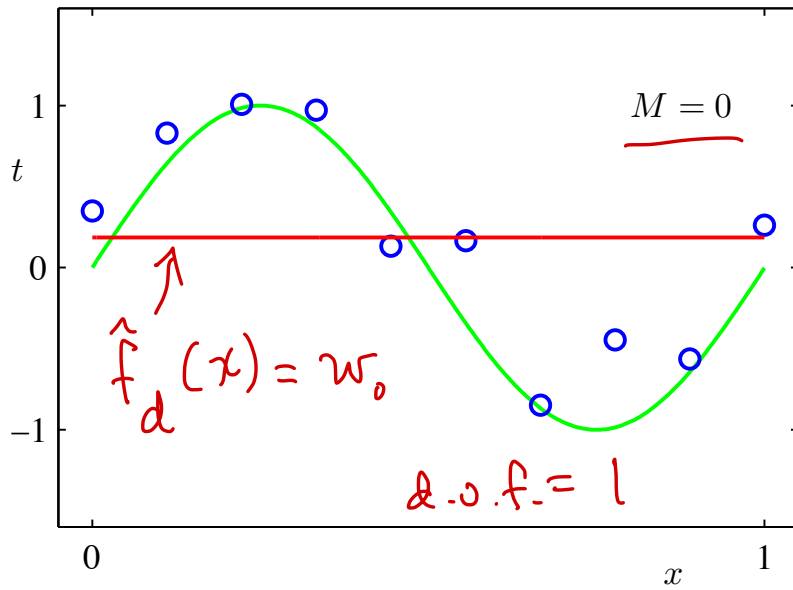
dim. 2 (augmented)

$N_c = ?$

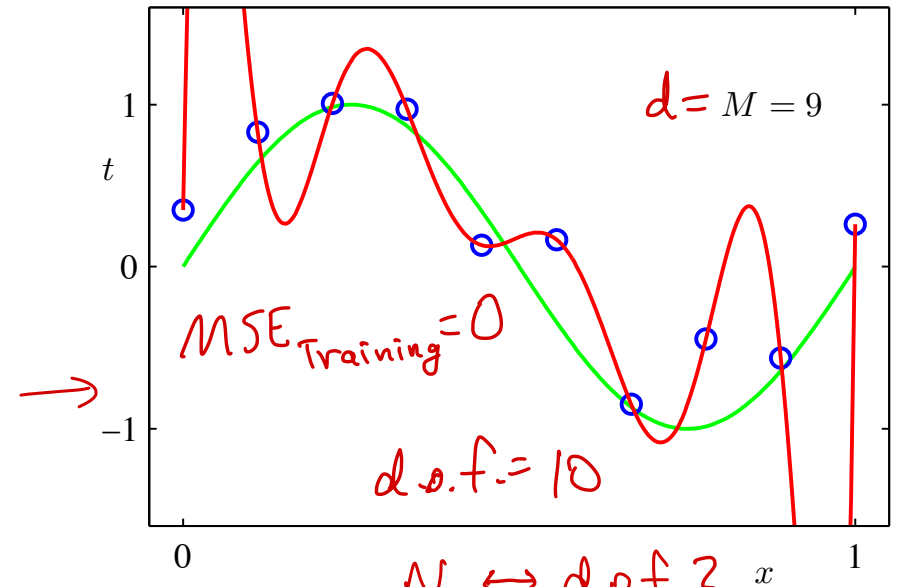
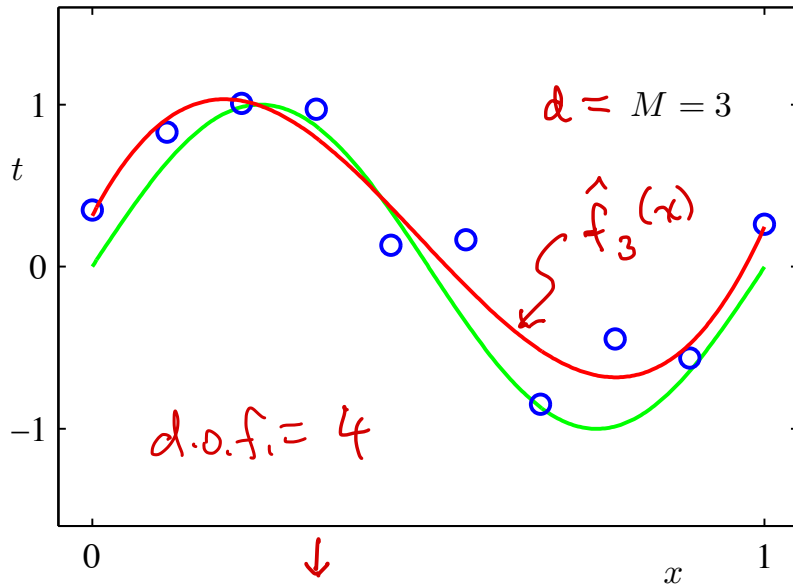
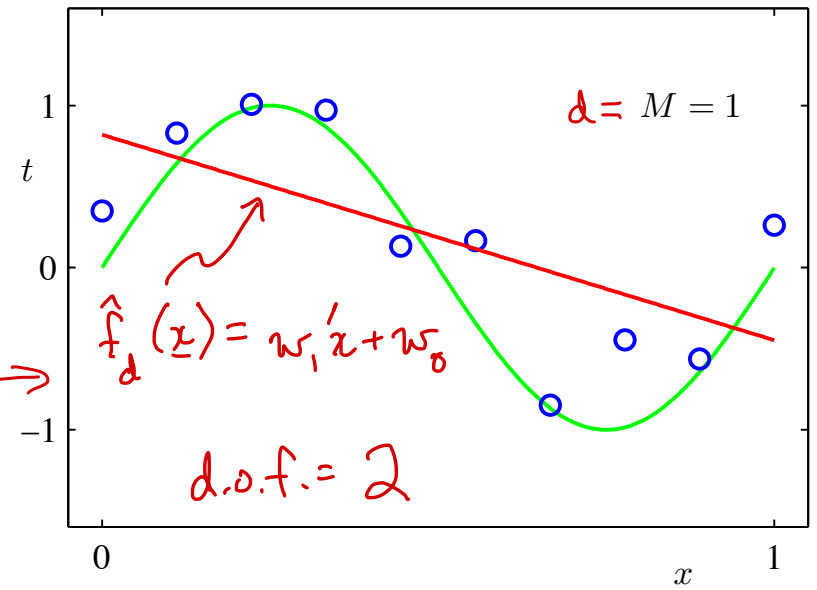
$N$  data points

Use MSE criterion  $\rightarrow$  least-squares regression.

# Regression complexity example: varying $d$



$N = 10$   
 $N_C = 10$   
 $d.o.f. = ?$



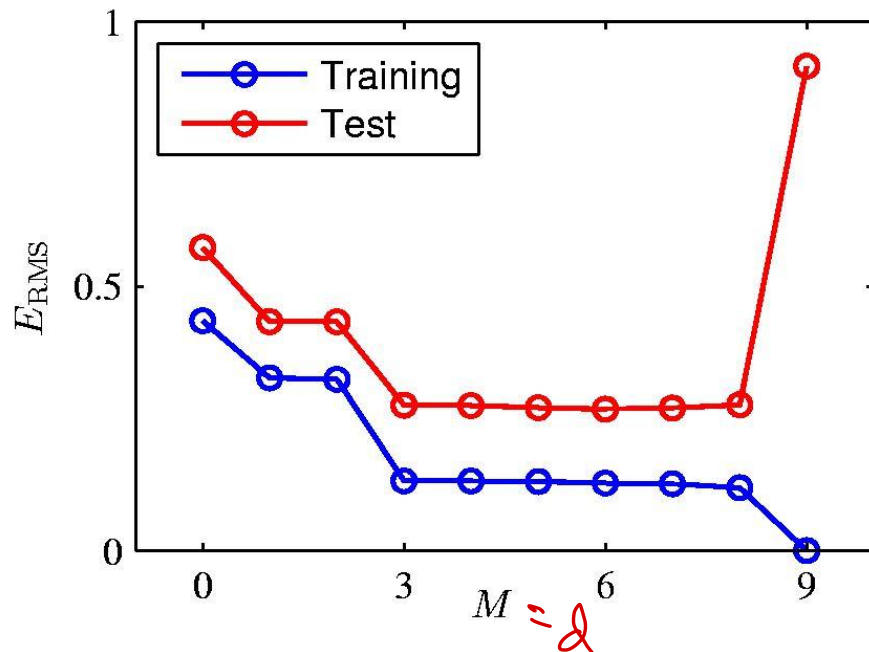
$N_C \leftrightarrow d.o.f.?$   
 $\downarrow$   
 $10 \quad \downarrow$   
 $4$

$N_C \leftrightarrow d.o.f.?$   
 $\downarrow$   
 $10 \quad \downarrow$   
 $10$

Regression complexity example: varying  $d$ 

$$N = 10$$

$$N_C = 10$$



## 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