

2019
怪兽
学堂



Regression

虾米

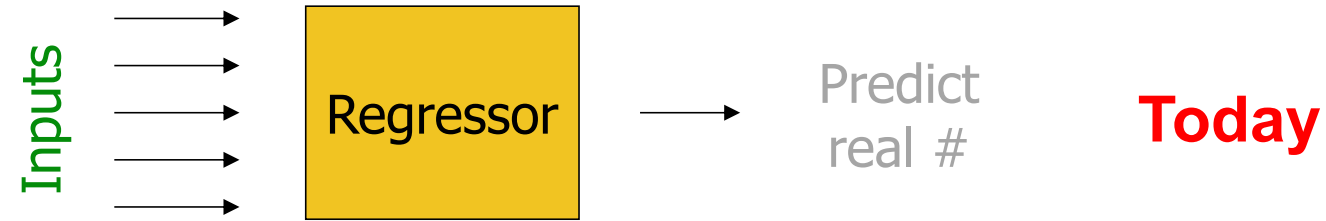
时间：2019年3月

Outline

- Regression
- Linear regression
 - As optimization → Gradient descent
- Overfitting and bias-variance

What is regression?

Where we are



Regression examples

Stock market



Weather prediction

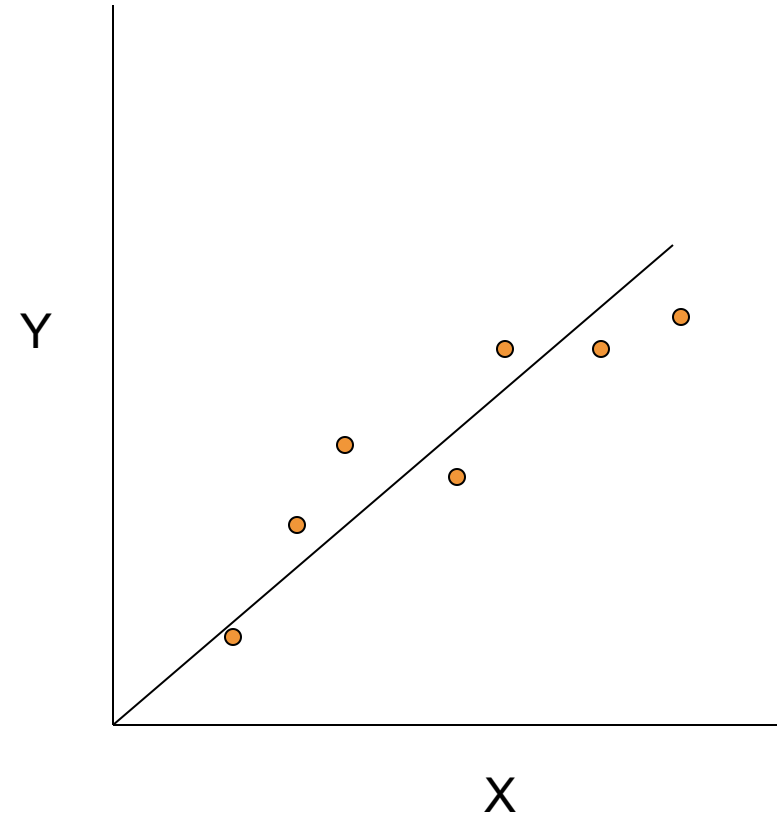


Temperature
72° F

Predict the temperature at any given location

Linear regression

- Given an input x we would like to compute an output y
- For example:
 - Predict height from age
 - Predict Google's price from Yahoo's price
 - Predict distance from wall from sensors

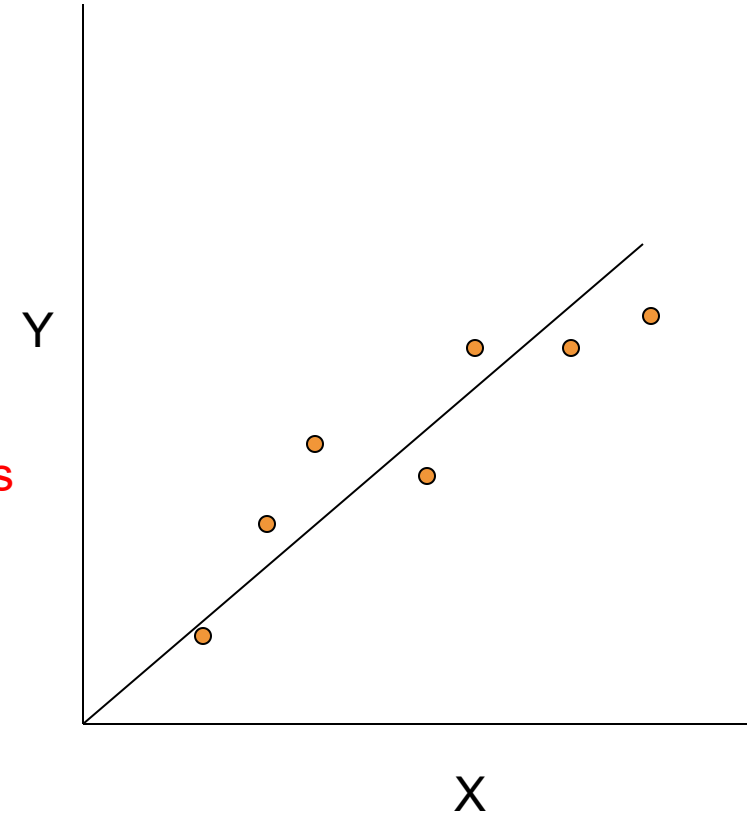


Linear regression

- Given an input x we would like to compute an output y
- In linear regression we assume that y and x are related with the following equation:

What we are trying to predict $y = wx + \varepsilon$ Observed values

where w is a parameter and ε represents measurement or other noise

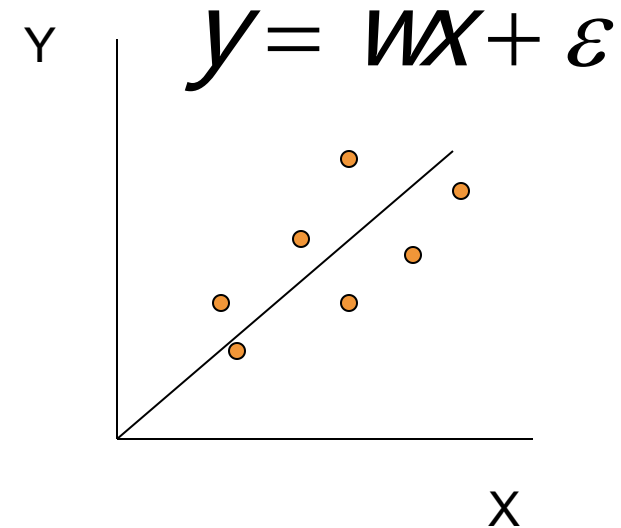


Linear regression

- Our goal is to estimate w from a training data of $\langle x_i, y_i \rangle$ pairs
- Optimization goal: minimize squared error (least squares):

$$\arg \min_w \sum_i (y_i - wx_i)^2$$

- Why least squares?
 - minimizes squared distance between measurements and predicted line
 - has a nice probabilistic interpretation
 - the math is pretty



Solving linear regression

- To optimize:
- We just take the derivative w.r.t. to w

$$\frac{\partial}{\partial w} \sum_i (y_i - wx_i)^2 = 2 \sum_i -x_i (y_i - wx_i)$$



prediction

Solving linear regression

- To optimize – closed form:
- We just take the derivative w.r.t. to w and set to 0:

$$\frac{\partial}{\partial w} \sum_i (y_i - wx_i)^2 = 2 \sum_i -x_i (y_i - wx_i) \Rightarrow$$

$$2 \sum_i x_i (y_i - wx_i) = 0 \Rightarrow 2 \sum_i x_i y_i - 2 \sum_i wx_i x_i = 0$$

$$\sum_i x_i y_i = \sum_i wx_i^2 \Rightarrow$$

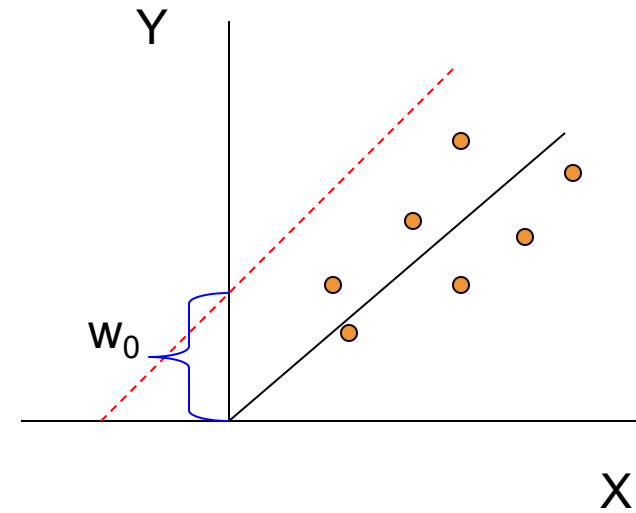
$$w = \frac{\sum_i x_i y_i}{\sum_i x_i^2}$$

Bias term

- So far we assumed that the line passes through the origin
- What if the line does not?
- No problem, simply change the model to

$$y = w_0 + w_1x + \varepsilon$$

- Can use least squares to determine w_0 , w_1



$$w_0 = \frac{\sum_i y_i - w_1 x_i}{n}$$

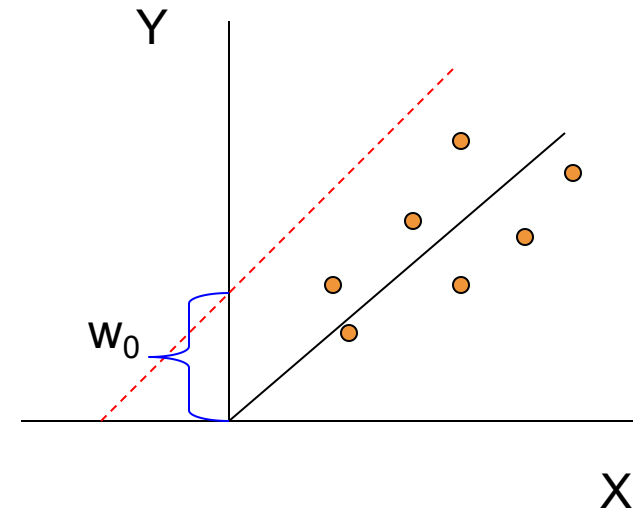
$$w_1 = \frac{\sum_i x_i (y_i - w_0)}{\sum_i x_i^2}$$

Bias term

- So far we assumed that the line passes through the origin
- What if the line does not?
- No problem, simply extend the model to

$$y = w_0 + w_1x + \varepsilon$$

- Can use least squares to determine w_0 , w_1

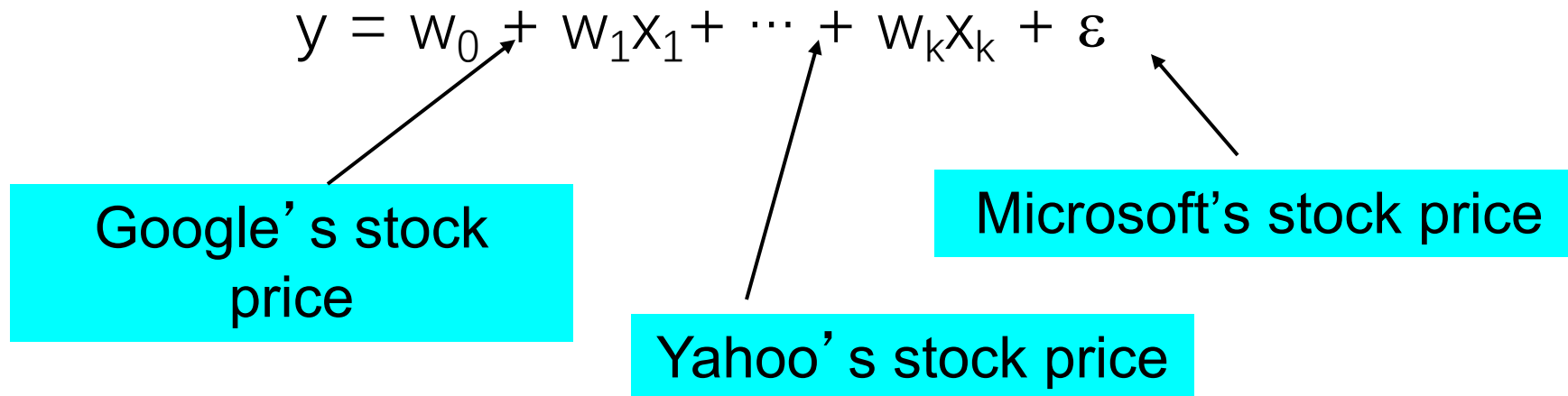


$$w_0 = \frac{\sum_i y_i - w_1 x_i}{n}$$

$$w_1 = \frac{\sum_i x_i (y_i - w_0)}{\sum_i x_i^2}$$

Multivariate regression

- What if we have several inputs?
 - Stock prices for Yahoo, Microsoft and Ebay for the Google prediction task
- This becomes a multivariate regression problem
- Again, its easy to model:

$$y = w_0 + w_1x_1 + \dots + w_kx_k + \varepsilon$$


The diagram illustrates the mapping of stock prices to the multivariate regression equation. Three cyan boxes at the bottom are labeled "Google's stock price", "Yahoo's stock price", and "Microsoft's stock price". Arrows point from these boxes to the equation above: "Google's stock price" points to the w_0 term, "Yahoo's stock price" points to the w_1x_1 term, and "Microsoft's stock price" points to the ε term.

Google's stock price

Yahoo's stock price

Microsoft's stock price

Multivariate regression

- What if we have several inputs?
 - Stock prices for Yahoo, Microsoft and Ebay for the Google prediction task
- This becomes a multivariate regression problem
- Again, its easy to model:

$$y = w_0 + w_1x_1 + \cdots + w_kx_k + \varepsilon$$

Not all functions can be approximated by a line/hyperplane...

$$y = 10 + 3x_1^2 - 2x_2^2 + \varepsilon$$

In some cases we would like to use polynomial or other terms based on the input data, are these still linear regression problems?

Non-Linear basis function

- So far we only used the observed values x_1, x_2, \dots
- However, linear regression can be applied in the same way to **functions** of these values
 - Eg: to add a term $w x_1 x_2$ add a new variable $z = x_1 x_2$ so each example becomes: x_1, x_2, \dots, z
- As long as these functions can be directly computed from the observed values the parameters are still linear in the data and the problem remains a multi-variate linear regression problem

$$y = w_0 + w_1 x_1^2 + \dots + w_k x_k^2 + \varepsilon$$

Non-Linear basis function

- How can we use this to add an intercept term?

Add a new “variable” $z=1$ and weight w_0

Non-linear basis functions

- What type of functions can we use?
- A few common examples:

- Polynomial: $\phi_j(x) = x^j$ for $j=0 \cdots n$

- Gaussian: $\phi_j(x) = \frac{(x - \mu_j)^2}{2\sigma_j^2}$

- Sigmoid: $\phi_j(x) = \frac{1}{1 + \exp(-s_j x)}$

- Logs: $\phi_j(x) = \log(x+1)$

Any function of the input values can be used. The solution for the parameters of the regression remains the same.

General linear regression problem

- Using our new notations for the basis function linear regression can be written as

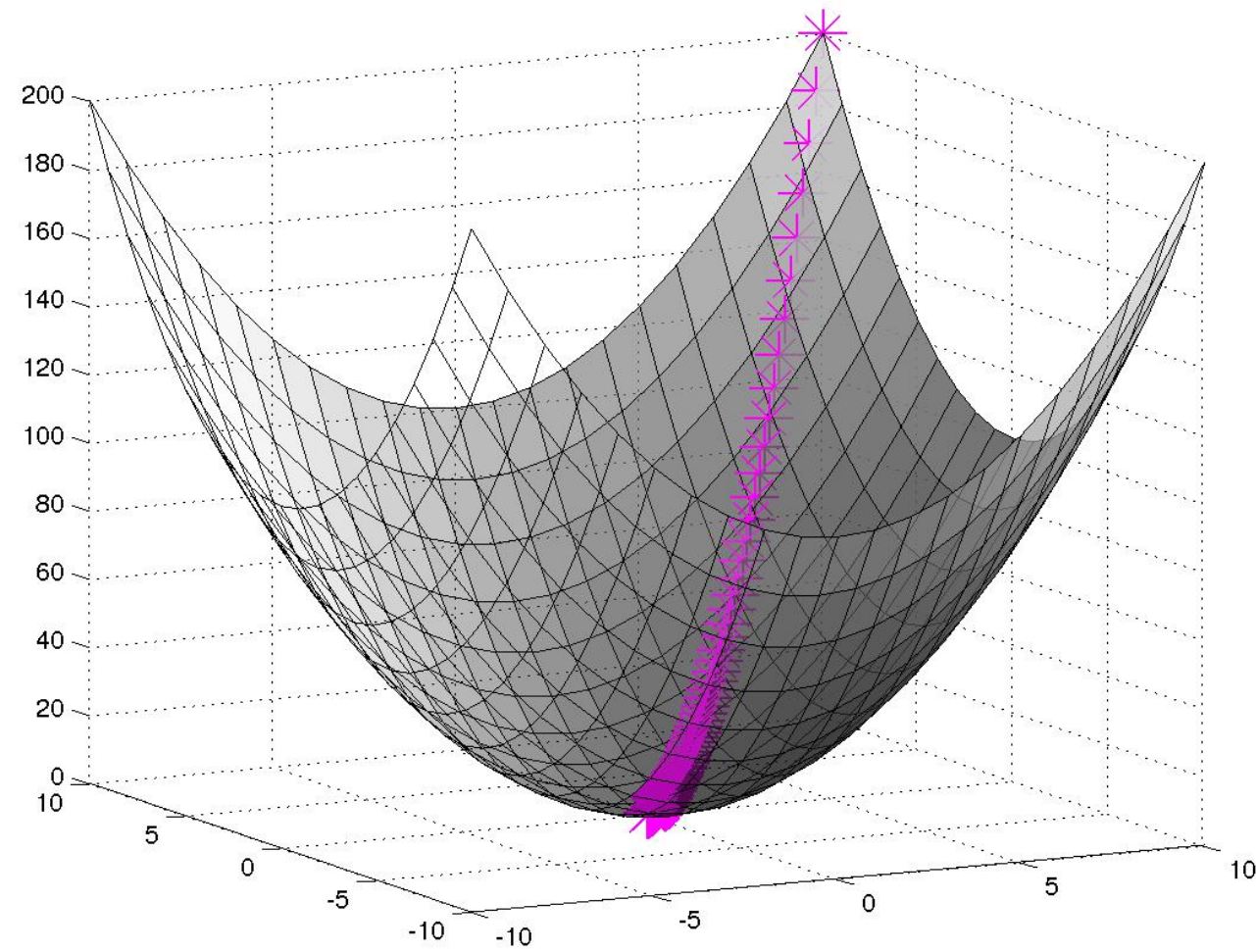
$$y = \sum_{j=0}^n w_j \phi_j(x)$$

- Where $\phi_j(\mathbf{x})$ can be either x_j for multivariate regression or one of the non-linear basis functions we defined
- ... and $\phi_0(\mathbf{x})=1$ for the intercept term

Learning/Optimizing Multivariate Least Squares

Approach 1: Gradient Descent

Gradient descent



Gradient Descent for Linear Regression

Goal: minimize the following loss function:

predict with: $\hat{y}^i = \sum_j^n w_j \phi_j(\mathbf{x}^i)$

$$J_{\mathbf{x},\mathbf{y}}(\mathbf{w}) = \sum_i (\mathcal{Y}^i - \hat{\mathcal{Y}}^i)^2 = \sum_i \left(\mathcal{Y}^i - \sum_j w_j \phi_j(\mathbf{x}^i) \right)^2$$

↑
sum over n examples

↑
sum over $k+1$ basis vectors

Gradient Descent for Linear Regression

Goal: minimize the following loss function:

predict with: $\hat{y}^i = \sum_j^n w_j \phi_j(\mathbf{x}^i)$

$$J_{\mathbf{x},\mathbf{y}}(\mathbf{w}) = \sum_i (\mathcal{Y}^i - \hat{\mathcal{Y}}^i)^2 = \sum_i \left(\mathcal{Y}^i - \sum_j w_j \phi_j(\mathbf{x}^i) \right)^2$$

$$\begin{aligned} \frac{\partial}{\partial w_j} J(\mathbf{w}) &= \frac{\partial}{\partial w_j} \sum_i (\mathcal{Y}^i - \hat{\mathcal{Y}}^i)^2 \\ &= 2 \sum_i (\mathcal{Y}^i - \hat{\mathcal{Y}}^i) \frac{\partial}{\partial w_j} \hat{\mathcal{Y}}^i \\ &= 2 \sum_i (\mathcal{Y}^i - \hat{\mathcal{Y}}^i) \frac{\partial}{\partial w_j} \sum_j w_j \phi_j(\mathbf{x}^i) \\ &= 2 \sum_i (\mathcal{Y}^i - \hat{\mathcal{Y}}^i) \phi_j(\mathbf{x}^i) \end{aligned}$$

Gradient Descent for Linear Regression

Learning algorithm:

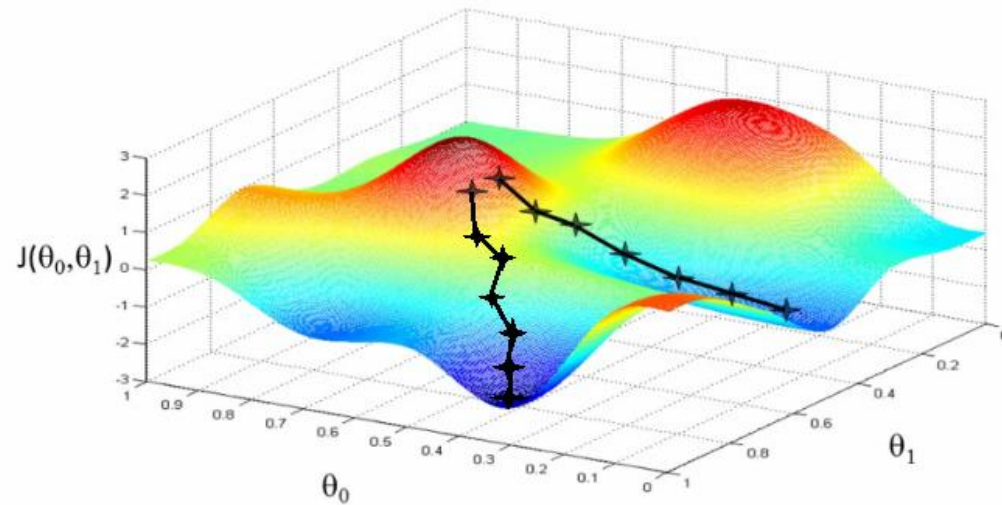
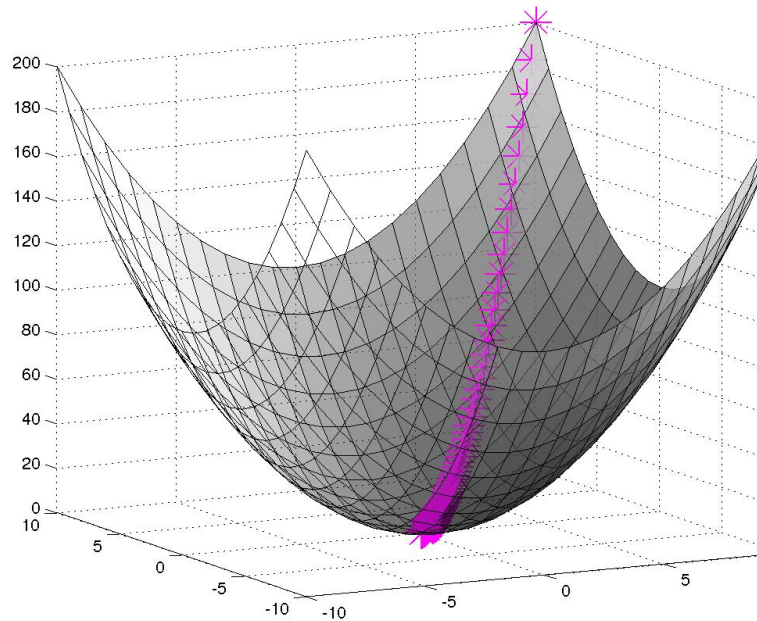
- Initialize weights **w=0**
- For t=1,... until convergence:
 - Predict for each example \mathbf{x}^i using **w**:
$$\hat{y}^i = \sum_{j=0}^k w_j \phi_j(\mathbf{x}^i)$$
 - Compute gradient of loss:
$$\frac{\partial}{\partial w_j} J(\mathbf{w}) = 2 \sum_i (y^i - \hat{y}^i) \phi_j(\mathbf{x}^i)$$
 - This is a vector **g**
 - Update: **w = w - λg**
 - λ is the learning rate.

Gradient Descent for Linear Regression

- We can use any of the tricks we used for logistic regression:
 - stochastic gradient descent (if the data is too big to put in memory)
 - regularization
 - ...

Linear regression is a *convex* optimization problem

so again gradient descent will reach a *global* optimum



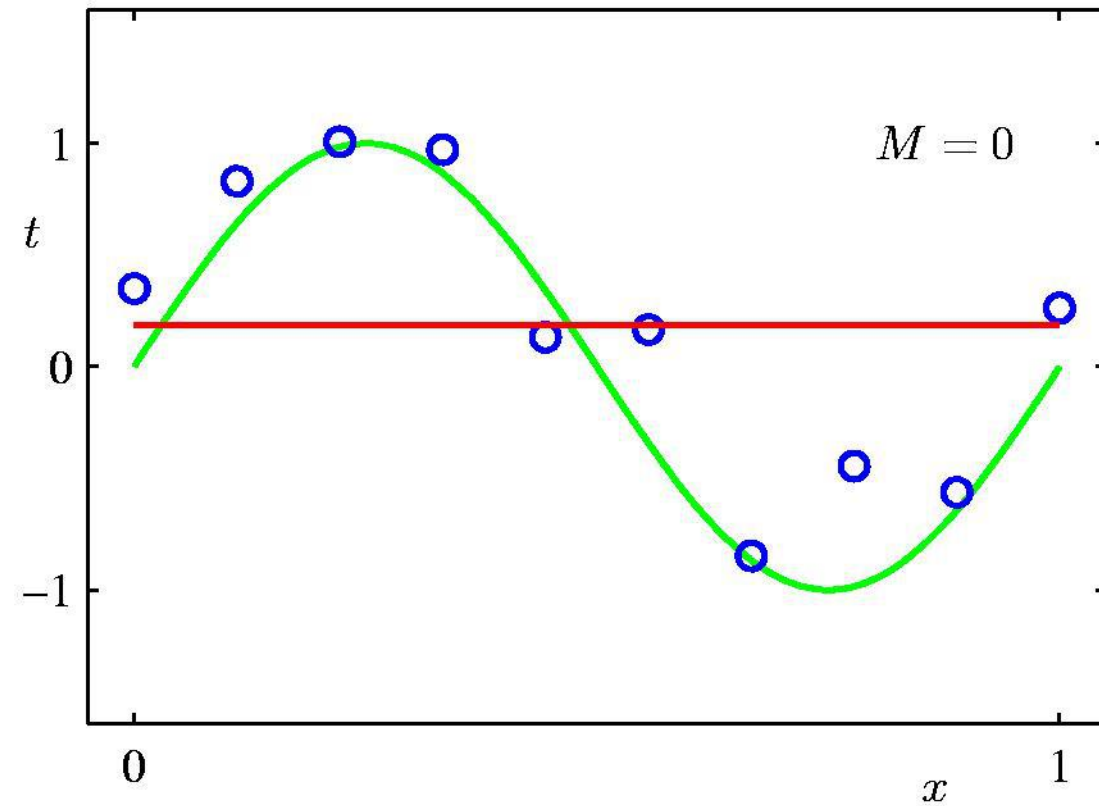
proof: differentiate again to get the second derivative

Regression and Overfitting

An example: polynomial basis vectors on a small dataset

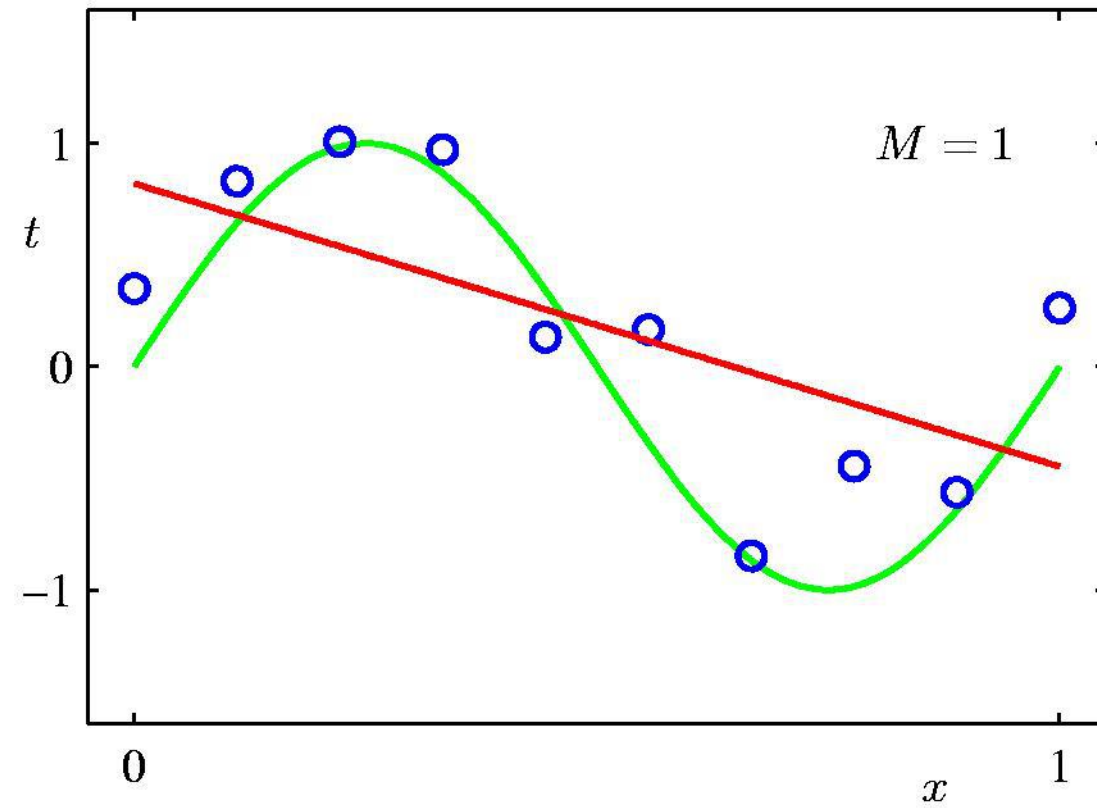
- From Bishop Ch 1

0th Order Polynomial

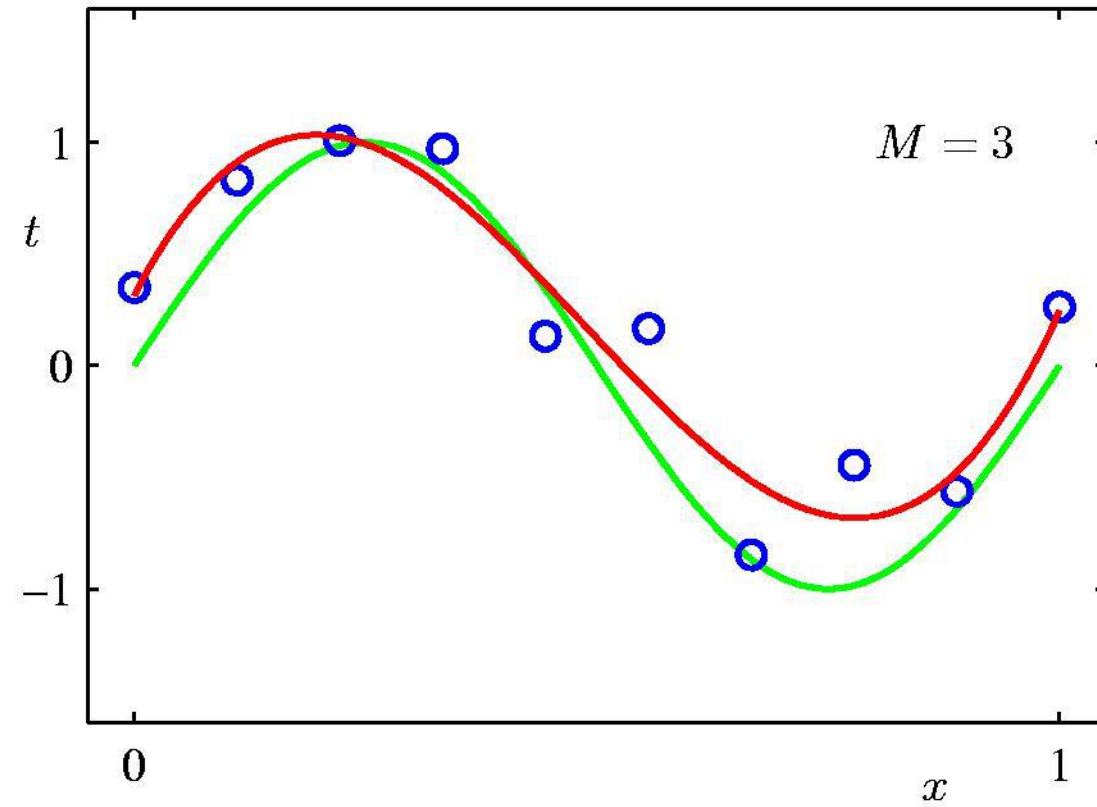


$n=10$

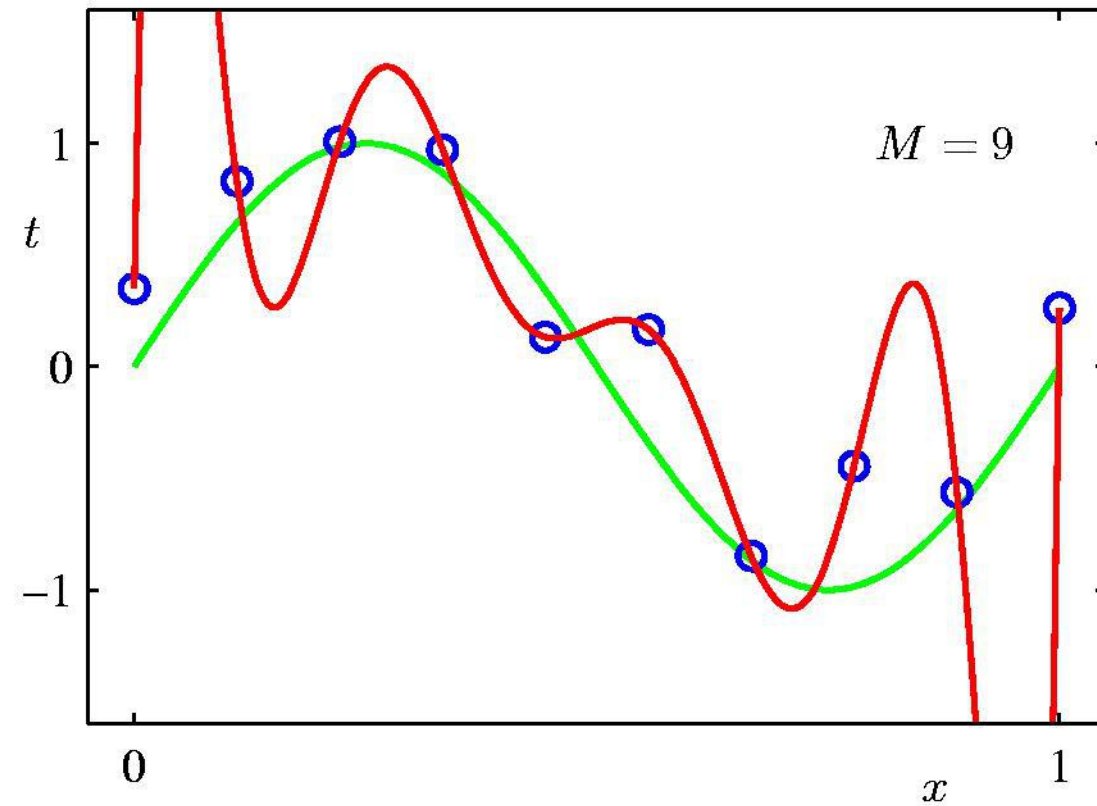
1st Order Polynomial



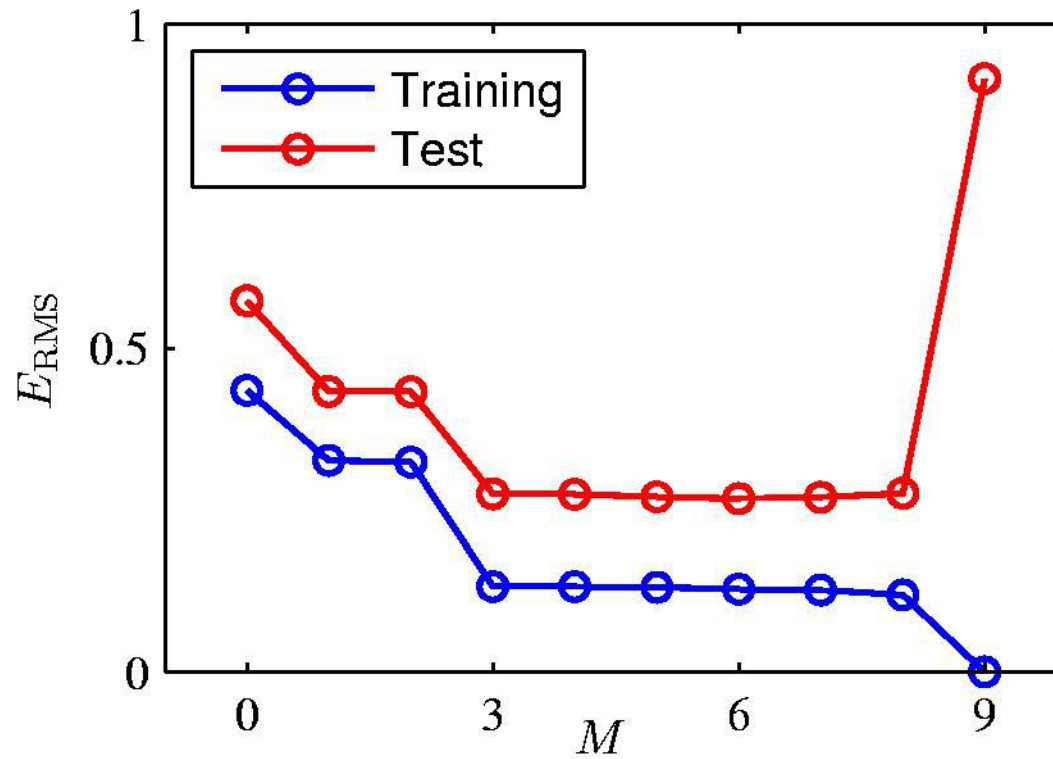
3rd Order Polynomial



9th Order Polynomial



Over-fitting



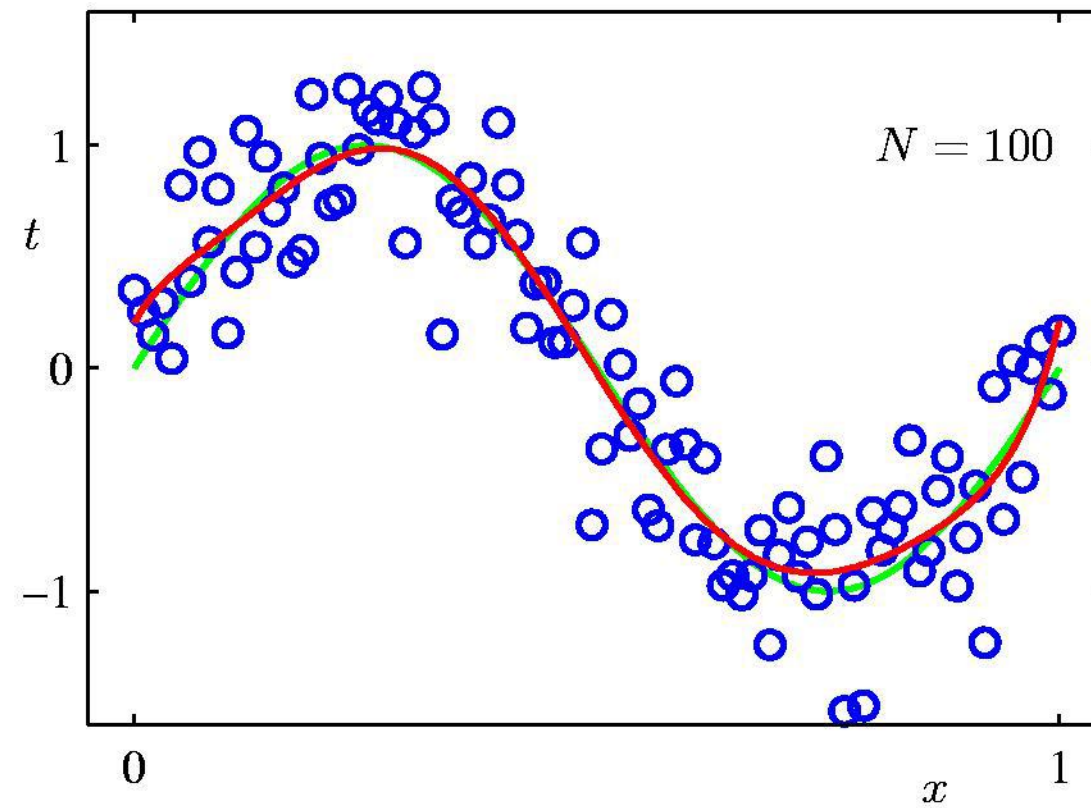
Root-Mean-Square (RMS) Error: $E_{\text{RMS}} = \sqrt{2E(\mathbf{w}^*)/N}$

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

Data Set Size:

9th Order Polynomial $N = 100$



Regularization

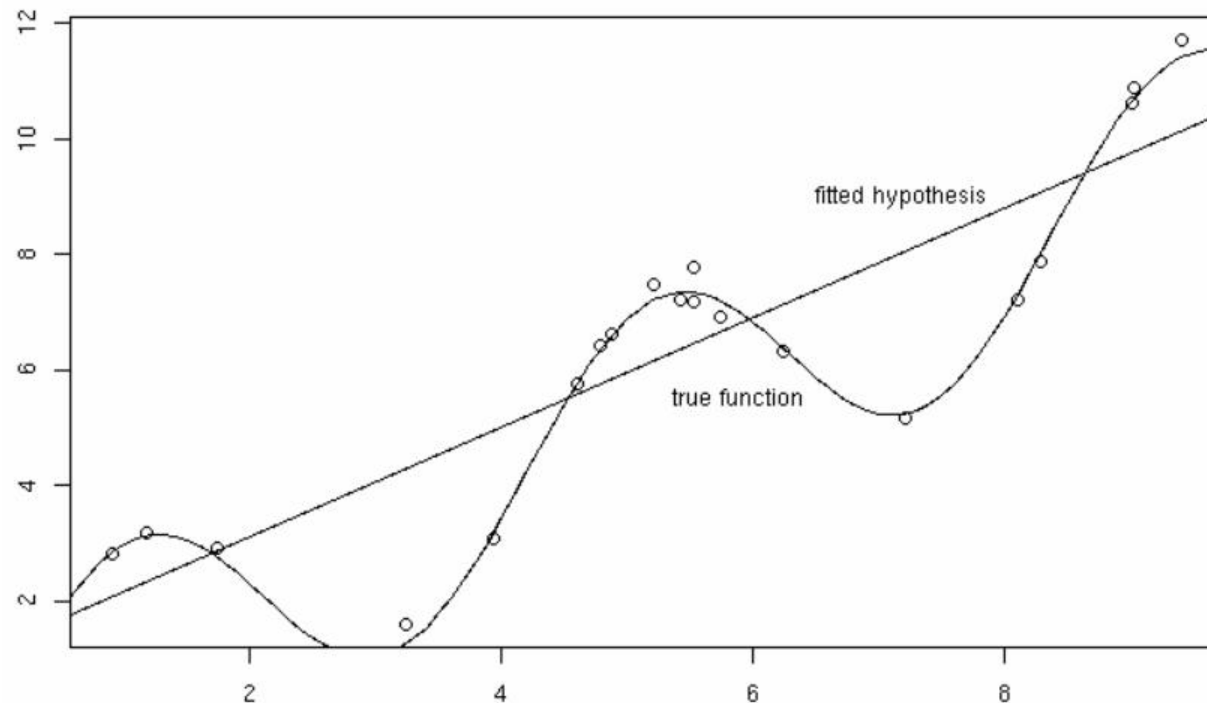
Penalize large coefficient values

$$J_{\mathbf{x},\mathbf{y}}(\mathbf{w}) = \frac{1}{2} \sum_i \left(y^i - \sum_j w_j \phi_j(\mathbf{x}^i) \right)^2 - \frac{\lambda}{2} \|\mathbf{w}\|^2$$

Understanding Overfitting: Bias-Variance

Example

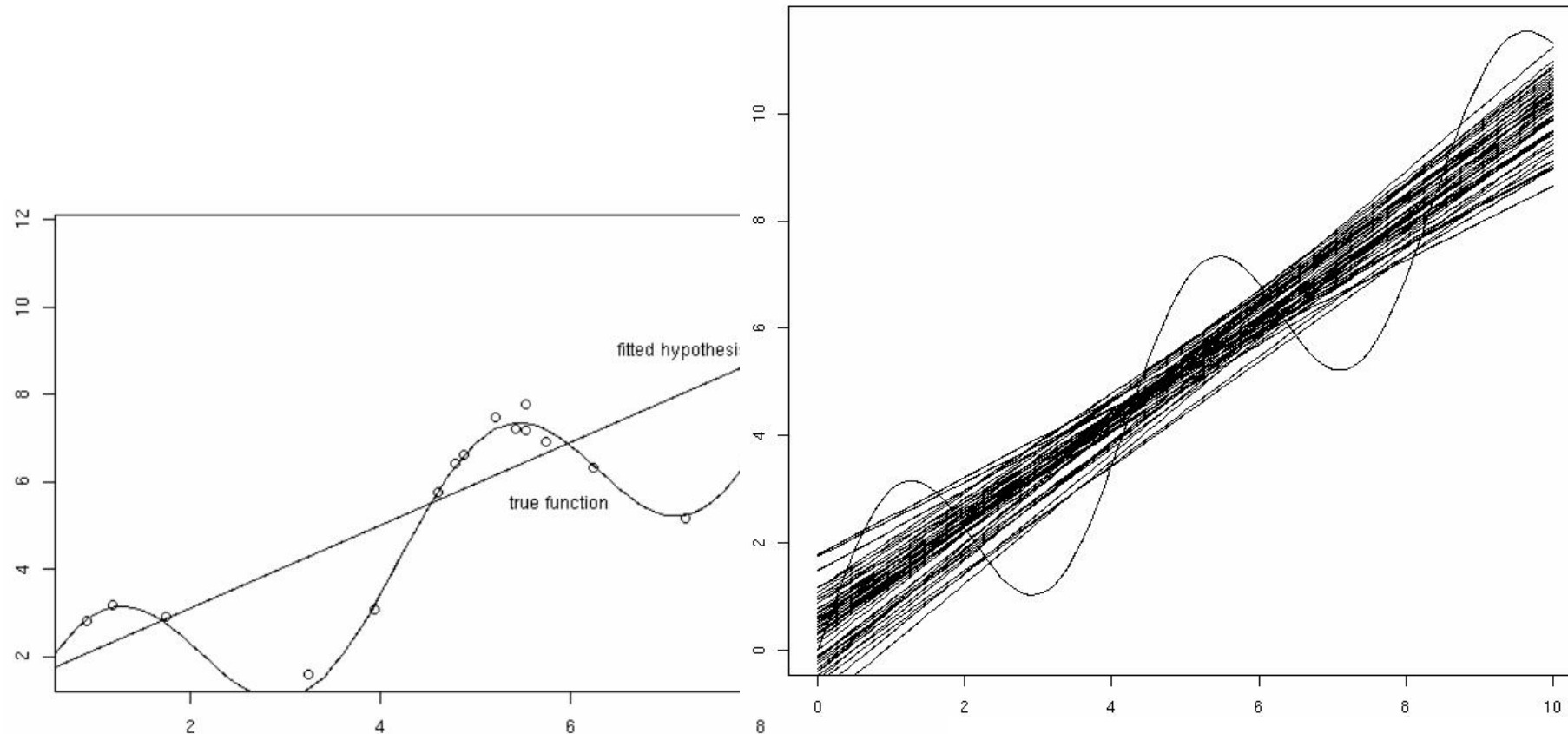
Tom Dietterich, Oregon St



$$y = x + 2 \sin(1.5x) + N(0,0.2)$$

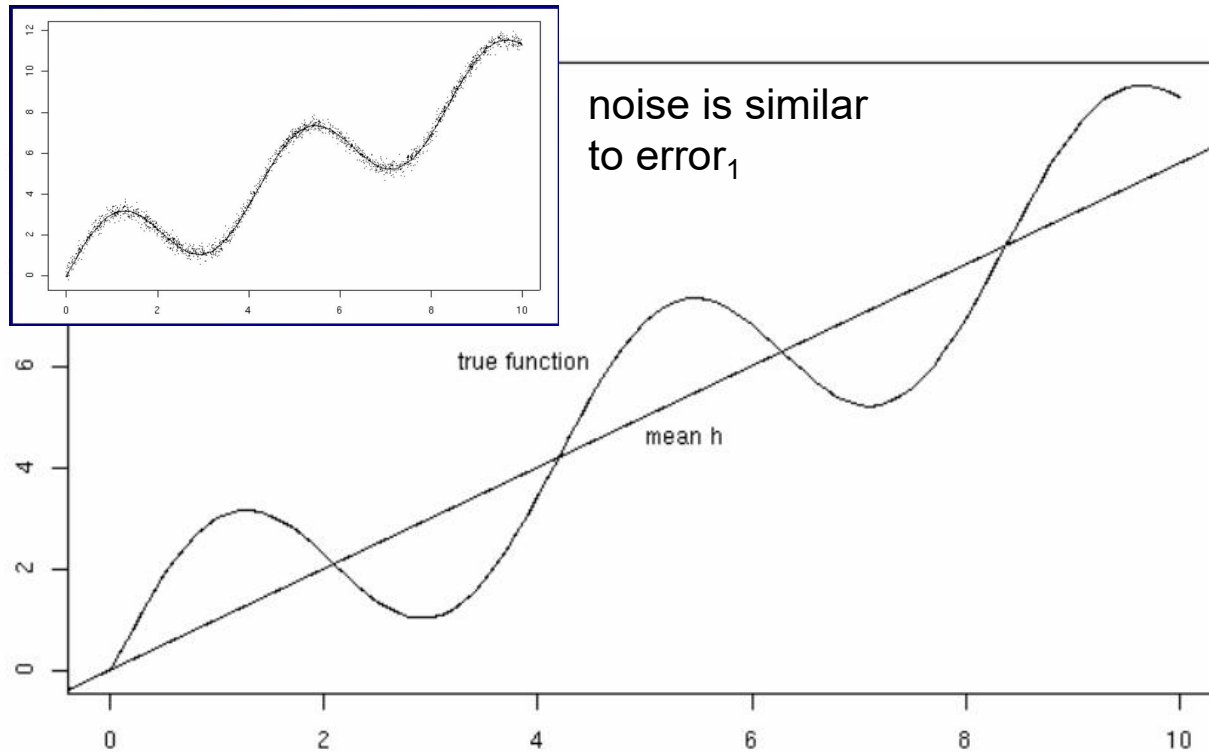
Example

Tom Dietterich, Oregon St



$$y = x + 2 \sin(1.5x) + N(0,0.2)$$

Same experiment, repeated:
with 50 samples of 20 points each



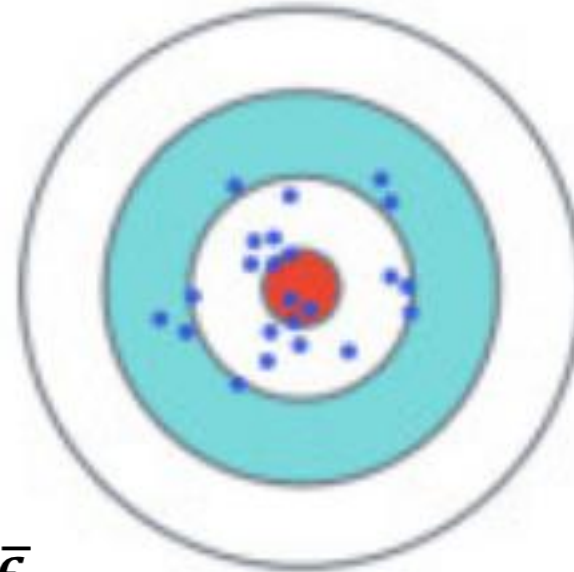
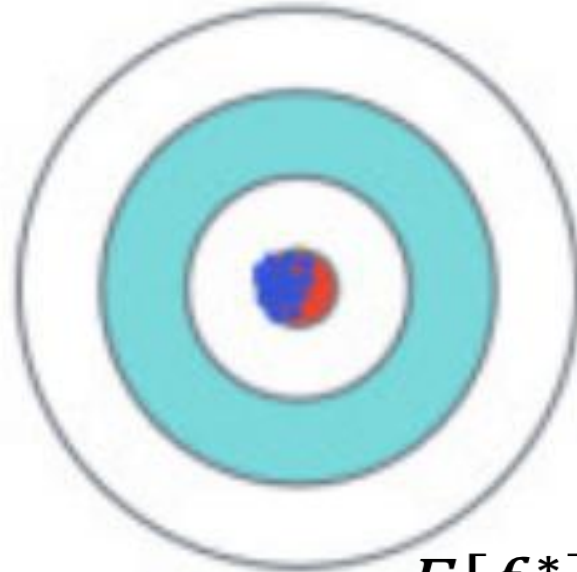
The true function f
can't be fit
perfectly with
hypotheses from
our class H
(lines) \rightarrow Error₁

Fix: *more*
expressive set of
hypotheses H

Low Variance

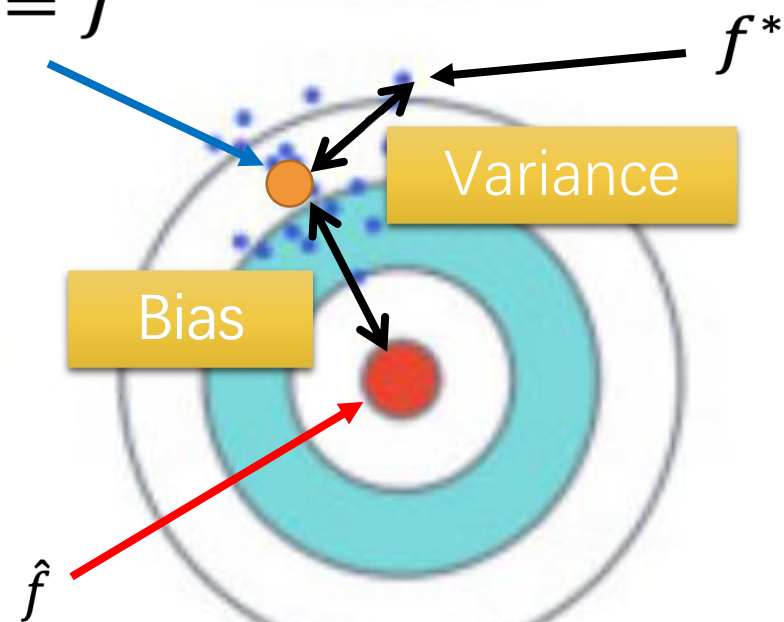
High Variance

Low Bias



$$E[f^*] = \bar{f}$$

High Bias

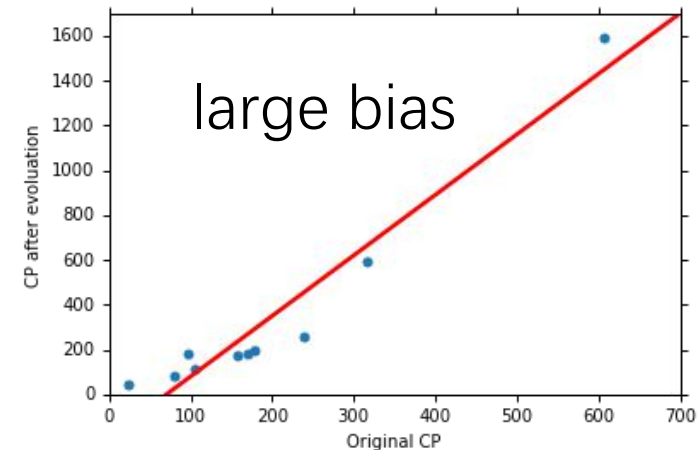


What to do with large bias?

- Diagnosis:
 - If your model cannot even fit the training examples, then you have large bias
- For bias, redesign your model:
 - Add more features as input
 - A more complex model

Underfitting

Overfitting



What to do with large variance?

- Diagnosis:
 - If you can fit the training data, but large error on testing data, then you probably have large variance
- For variance
 - More data: Very effective, but not always practical
 - Regularization

Overfitting

2019
怪兽
学堂



THANKS