

Machine Learning in Practice: a Crash Course

Lecture 4: Linear Regression & Bias-Variance

胡津铭
DolphinDB

Recap

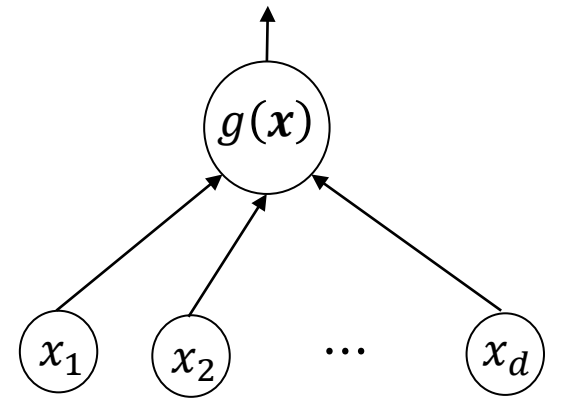
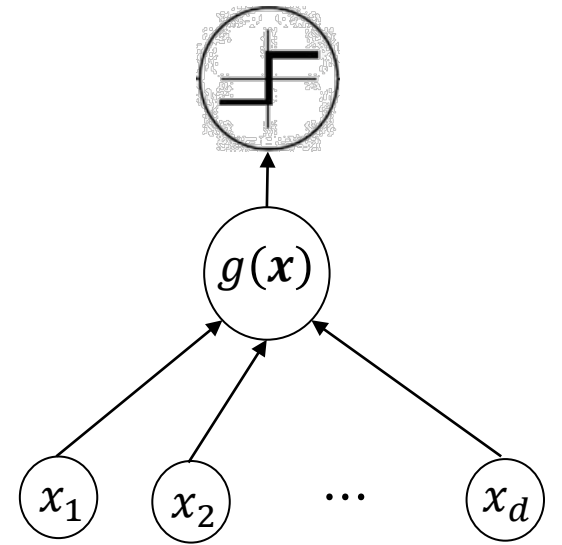
- Our goal (supervised learning):
 - To learn a set of discriminant functions
- Bayesian framework
 - We could design an optimal classifier if we knew:
 - $P(y_i)$: priors and $P(x | y_i)$: likelihood
 - Using training data to estimate $P(y_i)$ and $P(x | y_i)$
 - $P(y_i | x)$ is computed and be used as the discriminant functions
- Other possible ways?
 - Directly learning discriminant functions from the training data

Today

- Linear Regression
- Overfitting
- Bias-Variance

Classification VS Regression

- Both are supervised learning methods
 - Goal: learn a mapping from inputs \mathbf{x} to outputs y
- Classification (Categorization, Decision making...)
 - y is a categorical variable
- Regression
 - y is real-valued



Linear model

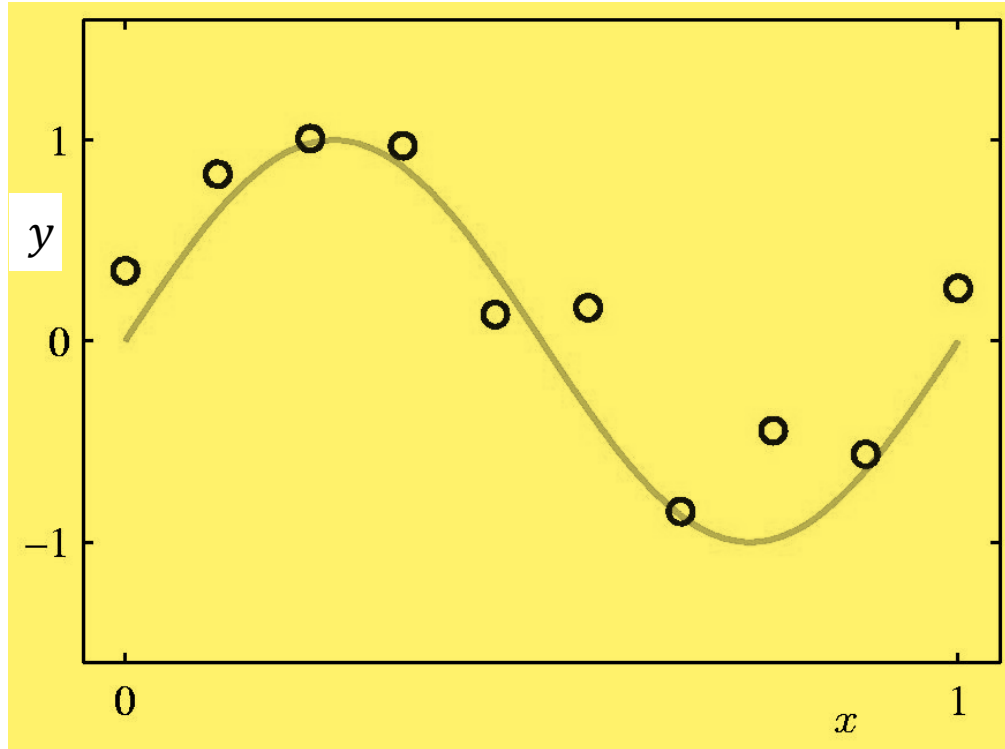
- Sample: $\mathbf{x} \in R^d, \mathbf{x} = [x_1, x_2, \dots x_d]^T$
- Finds a linear function $\mathbf{w} = [w_1, w_2, \dots, w_d]^T \in R^d, b$
- Representation: $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$$

$$\mathbf{x} = [x_1, x_2, \dots x_d, 1]^T \in R^{d+1}$$

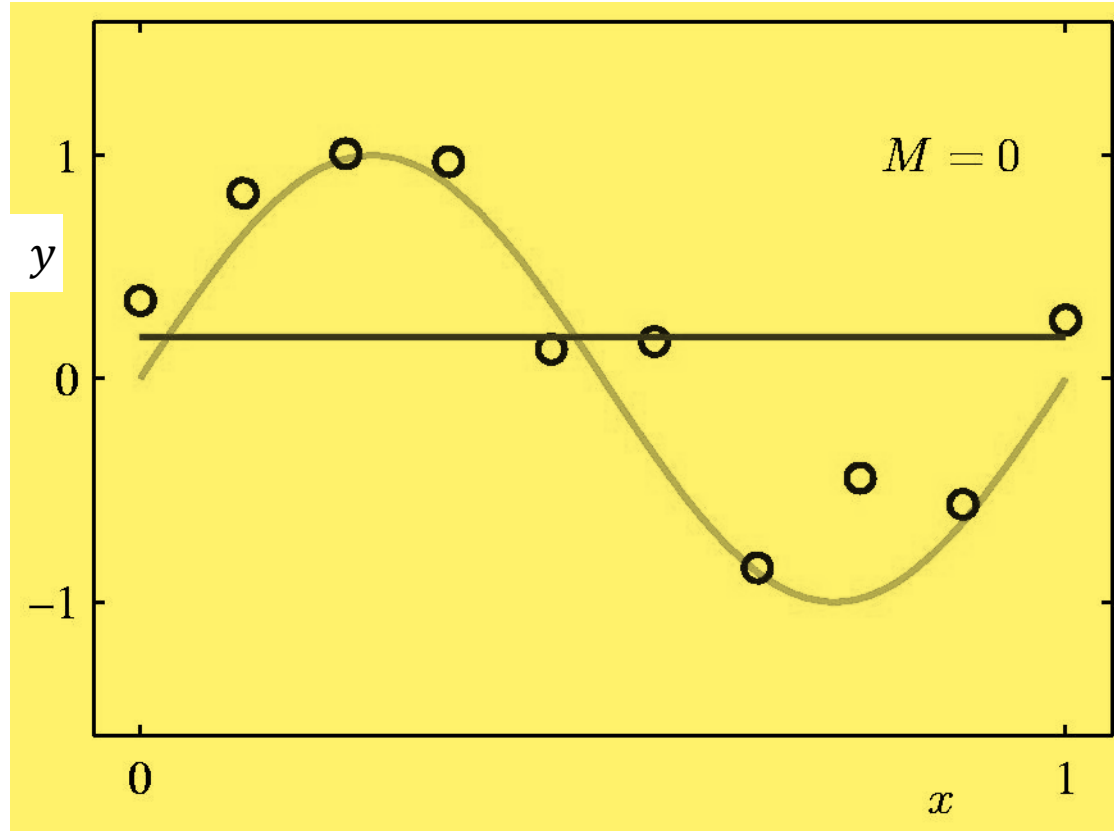
$$\mathbf{w} = [w_1, w_2, \dots w_d, b]^T \in R^{d+1}$$

Polynomial Curve Fitting



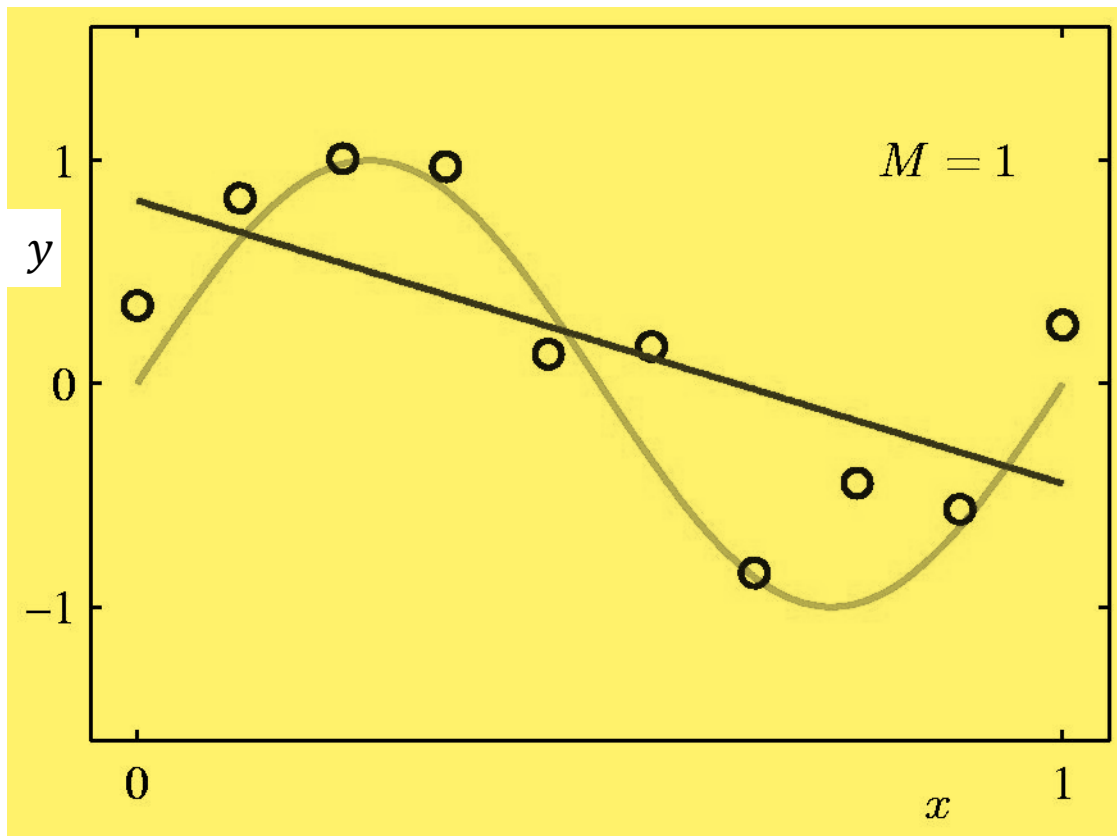
$$f(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \cdots + w_Mx^M = \sum_{j=0}^M w_jx^j$$

0th Order Polynomial



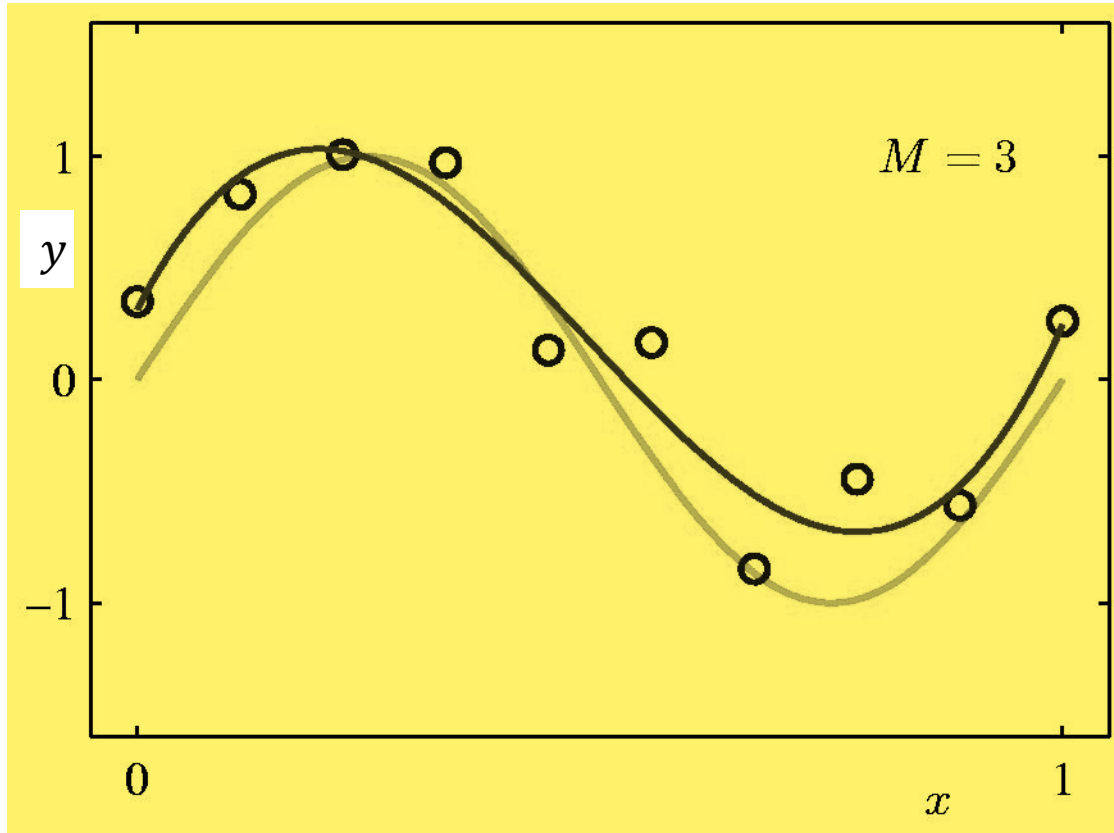
$$f(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \cdots + w_Mx^M = \sum_{j=0}^M w_jx^j$$

1st Order Polynomial



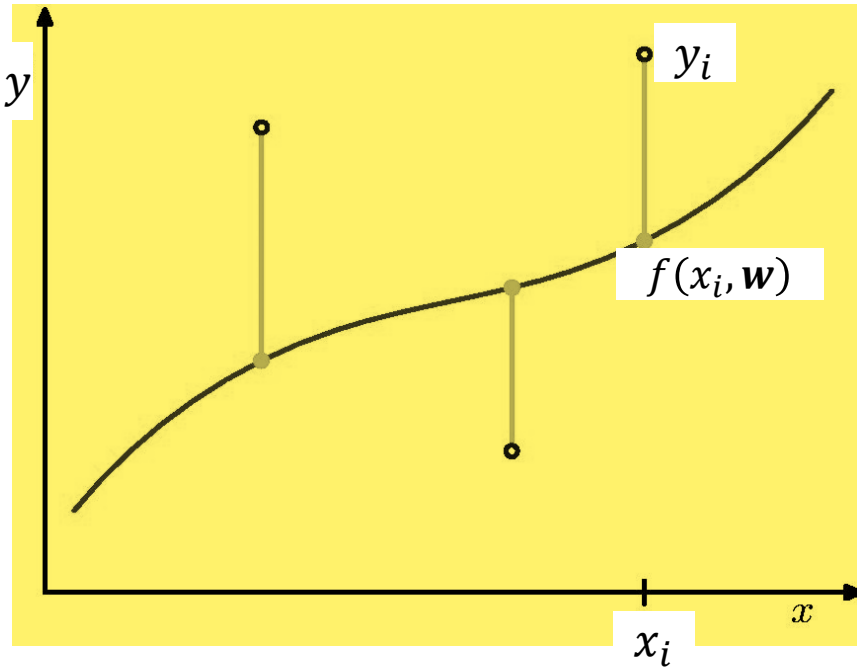
$$f(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \cdots + w_Mx^M = \sum_{j=0}^M w_jx^j$$

3rd Order Polynomial



$$f(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \cdots + w_Mx^M = \sum_{j=0}^M w_jx^j$$

Sum-of-Squares Error Function



- Training data:
 - $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$
- To learn f which $f(x) = y$
- Evaluation & Loss:
$$\text{MSE}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i, \mathbf{w}))^2$$

Linear Regression Model

- Training data: (\mathbf{x}_i, y_i)
- $f(\mathbf{x}) = w_0 + \sum_{i=1}^d w_i x_i = w_0 + \mathbf{w}^T \mathbf{x}$
 - $\mathbf{w} = [w_1, \dots, w_d]^T$ and w_0 : unknown parameters or coefficients
 - \mathbf{x} : Feature vector, the outcome of **feature engineering/extraction**.
- Minimize the **mean-squared error** :

$$J = \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2$$

How to learn the parameters?

- Optimize against the loss function!

$$J = \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2$$

- How to?
 1. Least Square Error Solver
 2. Gradient Descent

Least Square Errors Solver

$$J = \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2$$

- Using matrix notation for convenience

$$X = [\mathbf{x}_1, \dots, \mathbf{x}_n], \quad \mathbf{y} = [y_1, \dots, y_n]^T$$

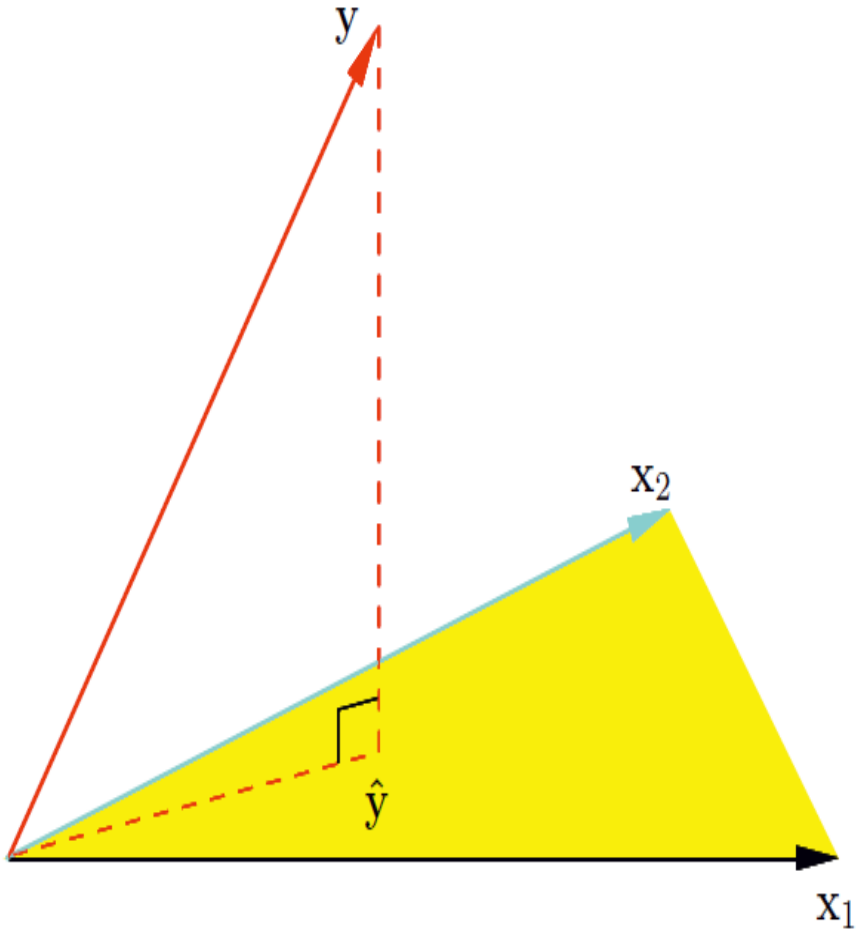
- Loss(\mathbf{w}) = $(\mathbf{y} - X^T \mathbf{w})^T (\mathbf{y} - X^T \mathbf{w})$
- Computing the gradient/derivative gives:

$$\nabla J = -2X(\mathbf{y} - X^T \mathbf{w})$$

- Setting the gradient to zero

$$\begin{aligned} XX^T \mathbf{w} &= X\mathbf{y} \\ \mathbf{w} &= (XX^T)^{-1} X\mathbf{y} \end{aligned}$$

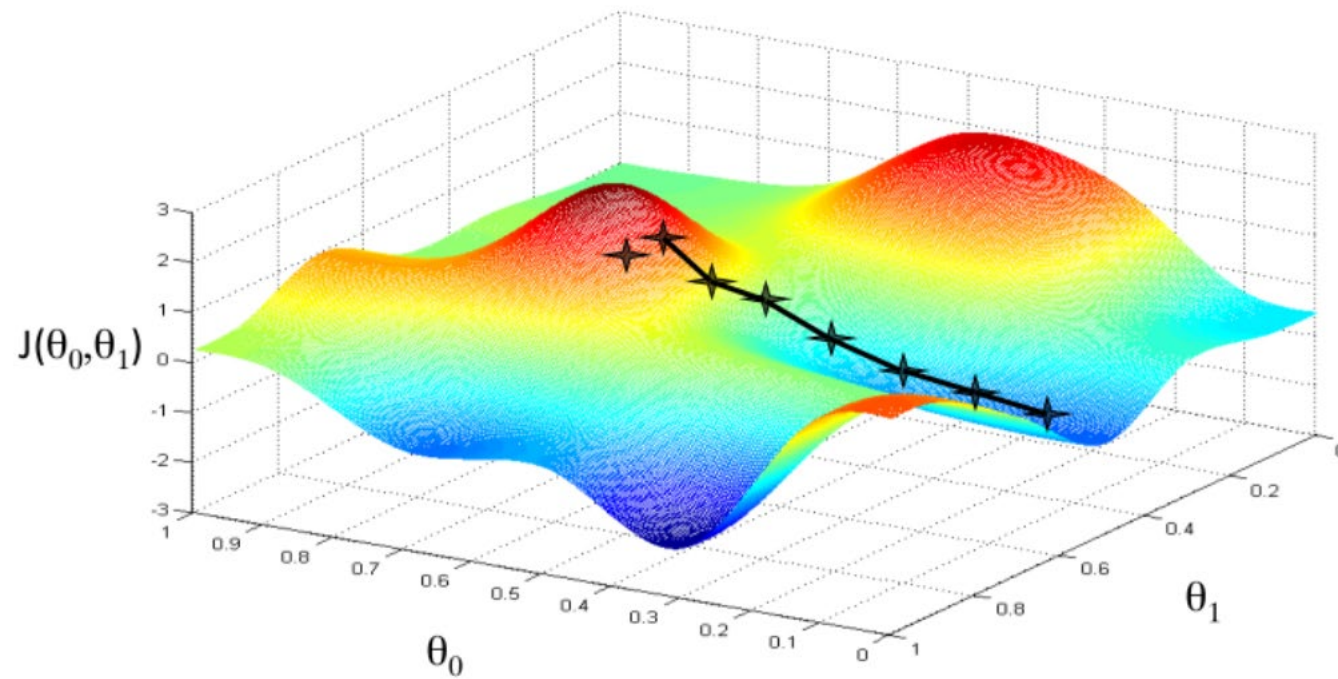
Geometry of least-squares fitting



Gradient Descent



Gradient Descent

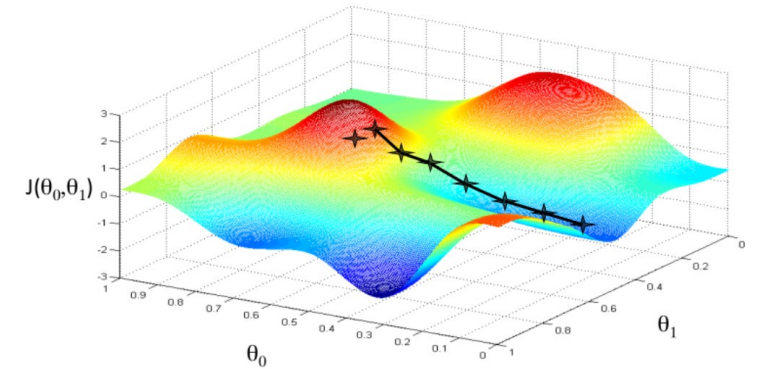


Gradient Descent

- It can find a local minimum of a function
- One takes steps proportional to the negative of the gradient (or of the approximate gradient) of the function at the current point.

Algorithm 1 (Basic gradient descent)

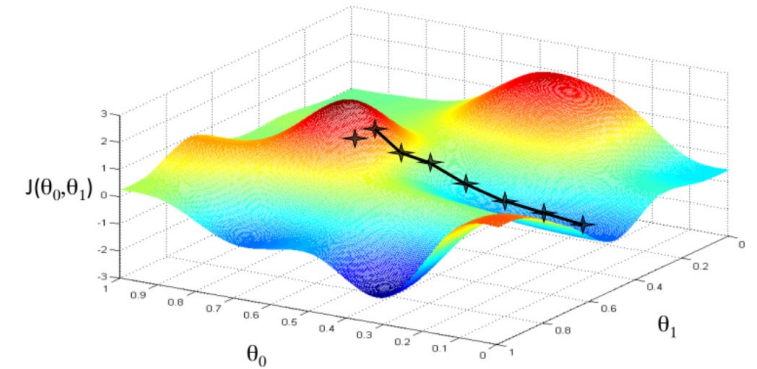
```
1 begin initialize  $\mathbf{a}$ , criterion  $\theta, \eta(\cdot), k = 0$   
2   do  $k \leftarrow k + 1$   
3      $\mathbf{a} \leftarrow \mathbf{a} - \eta(k) \nabla J(\mathbf{a})$   
4   until  $\eta(k) \nabla J(\mathbf{a}) < \theta$   
5 return  $\mathbf{a}$   
6 end
```



Gradient Descent

Algorithm 1 (Basic gradient descent)

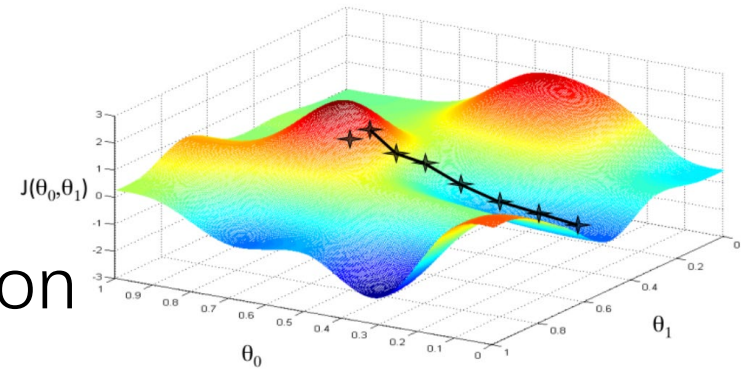
```
1 begin initialize  $\mathbf{a}$ , criterion  $\theta, \eta(\cdot), k = 0$   
2   do  $k \leftarrow k + 1$   
3      $\mathbf{a} \leftarrow \mathbf{a} - \eta(k) \nabla J(\mathbf{a})$   
4   until  $\eta(k) \nabla J(\mathbf{a}) < \theta$   
5 return  $\mathbf{a}$   
6 end
```



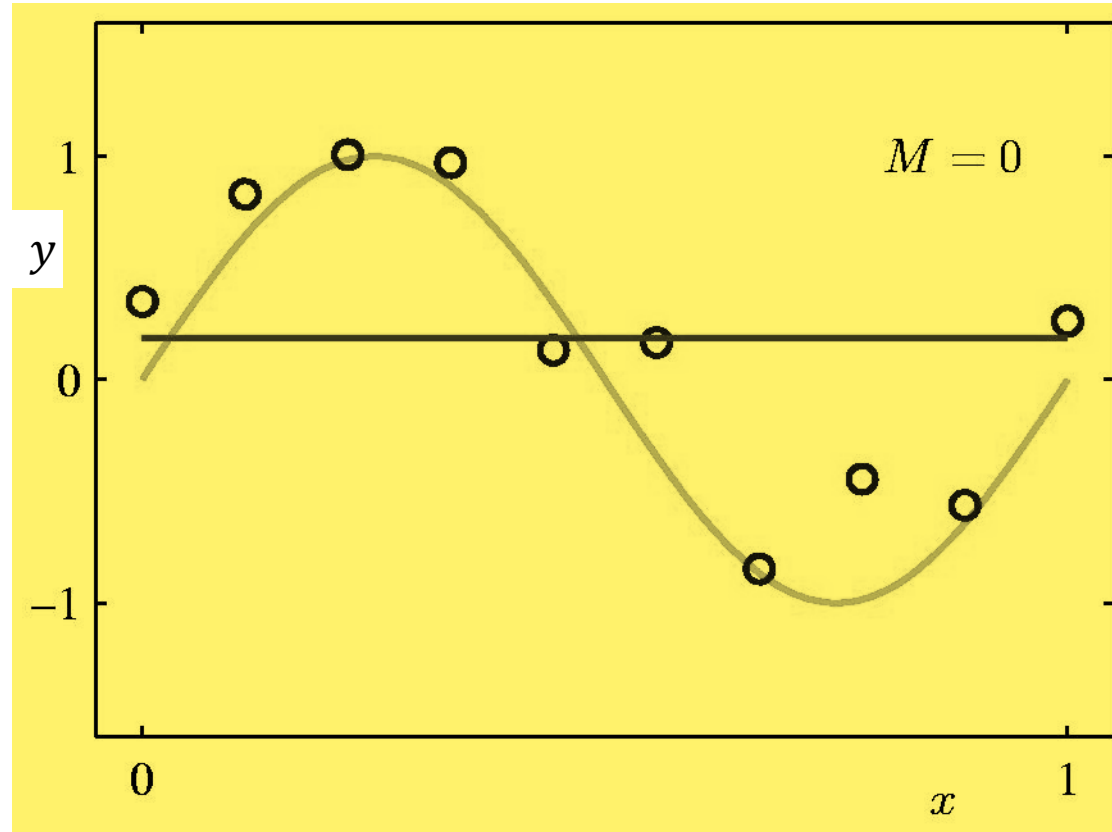
- Time complexity: $O(ndt)$
- Stochastic Gradient Descent(SGD)
 - Randomly select b (often called batch size) samples from the data
 - Time complexity: $O(bdt)$
- SGD(and its variants) are widely used, especially in deep learning

Convex Optimization

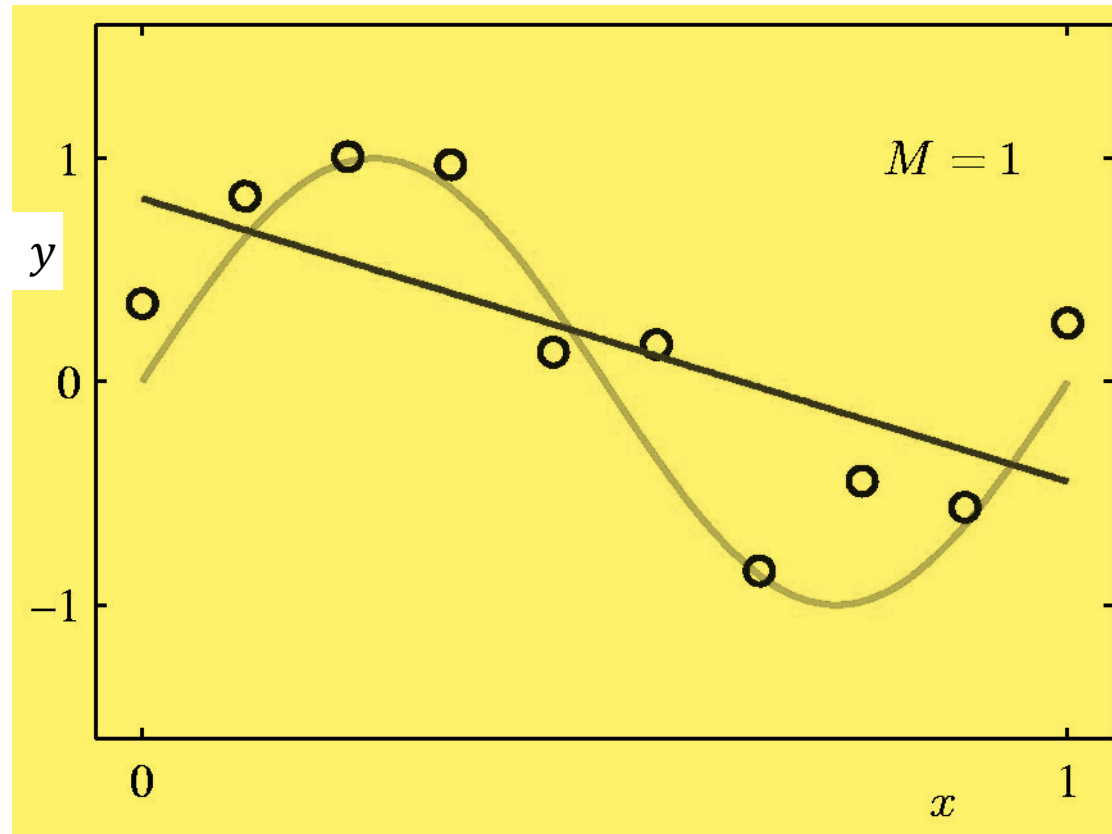
- Least Square Error is a so-called convex problem
- There are a bunch of optimization methods for such problems, called convex optimization
- Gradient Descent: first-order
- Newton's method: second-order
- Quasi Newton's methods
-
- If the problem is nonconvex, then can only find the local minimum/optima



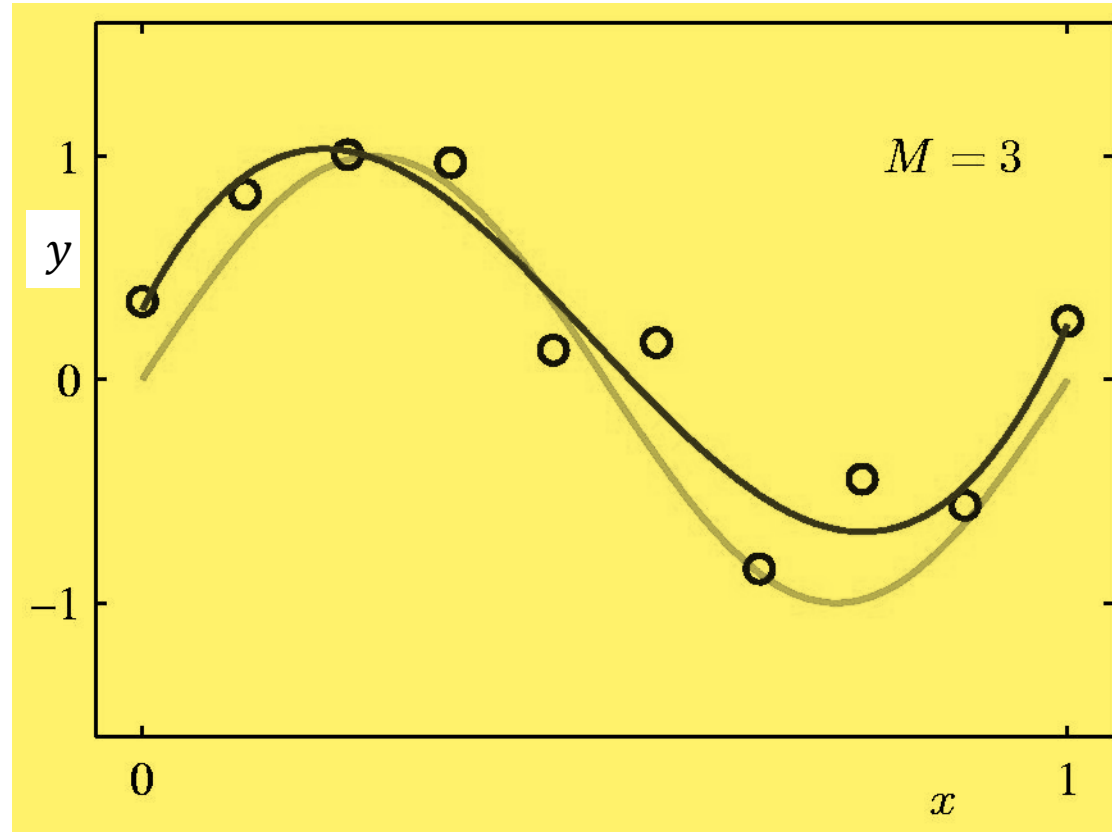
0th Order Polynomial



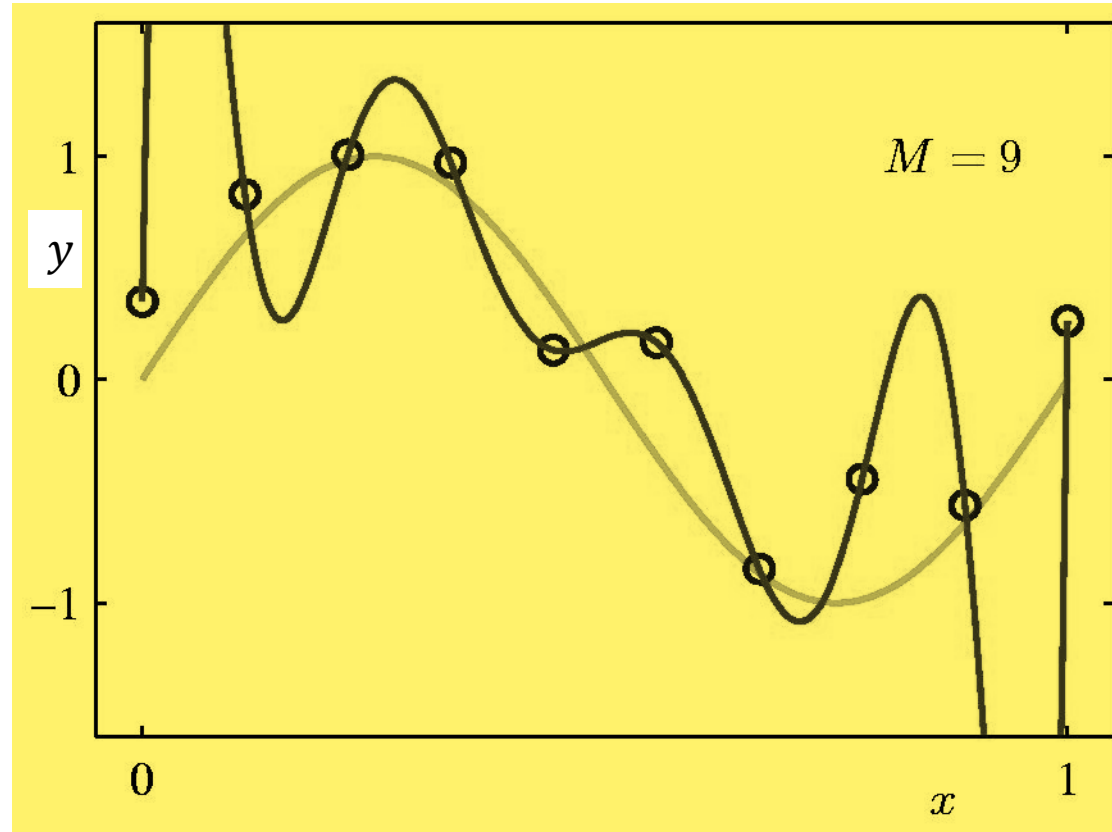
1st Order Polynomial



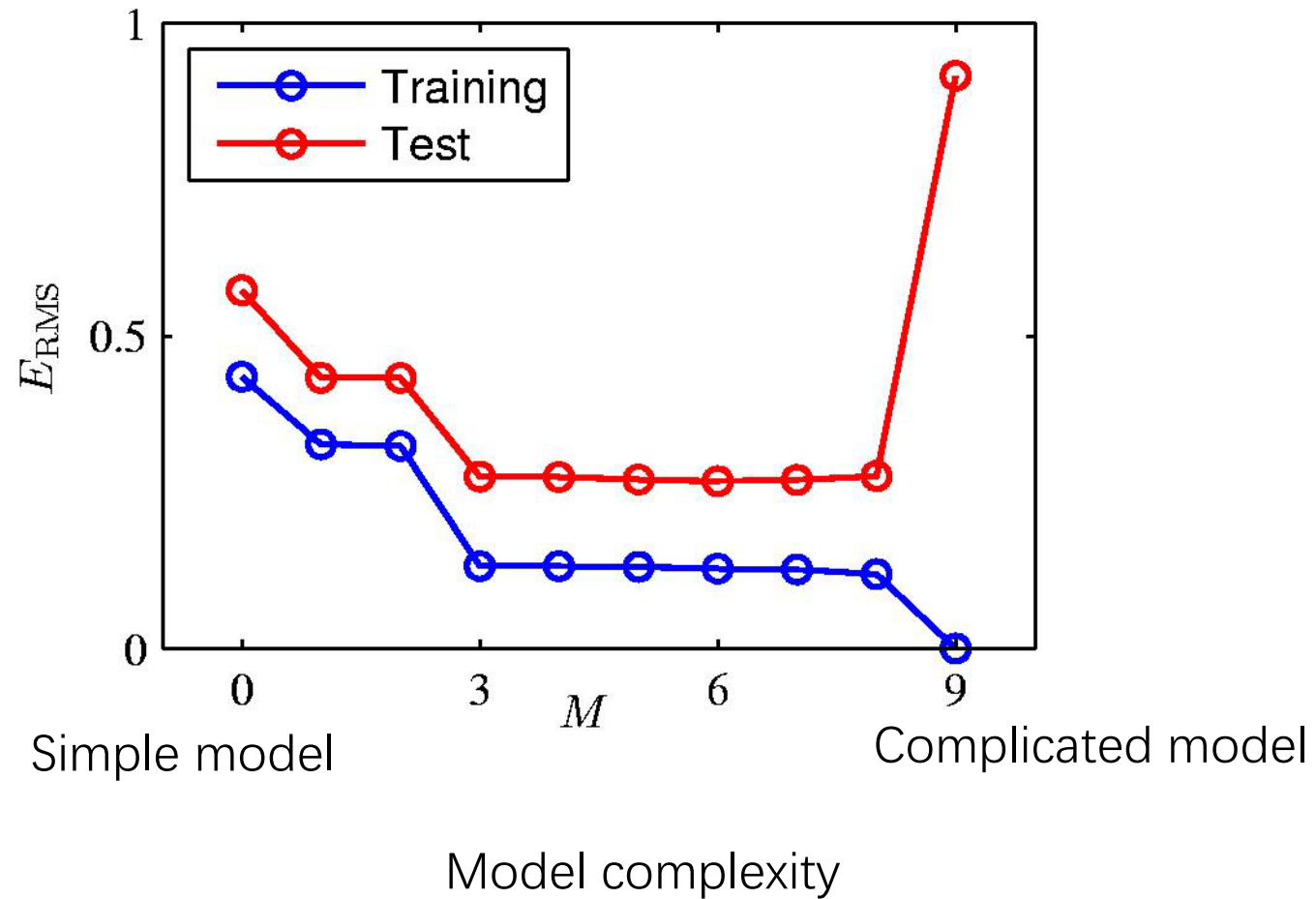
3rd Order Polynomial



9th Order Polynomial

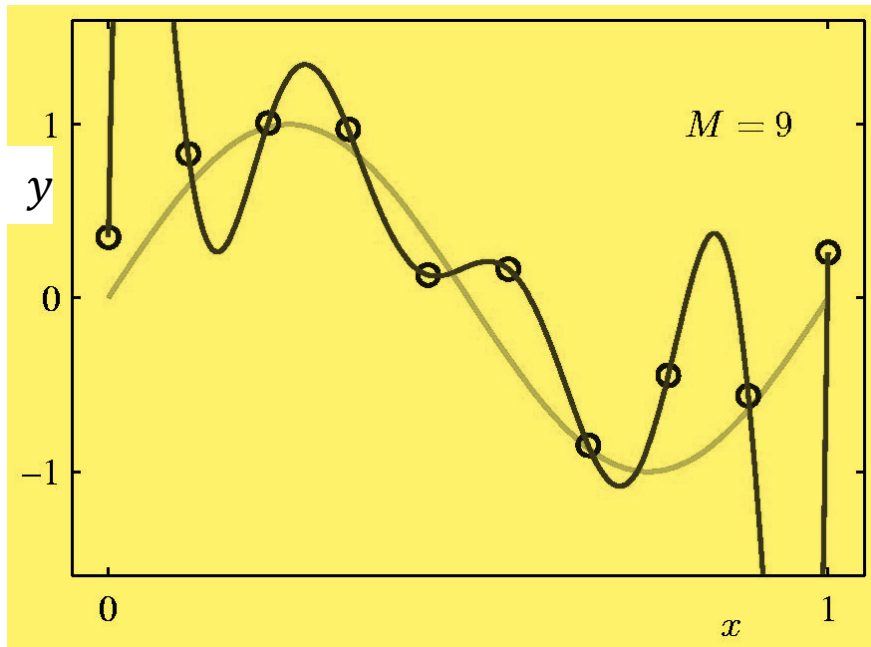


Overfitting



Issues with MSE Criterion

- Coefficients
- Overfit the noise in the data!



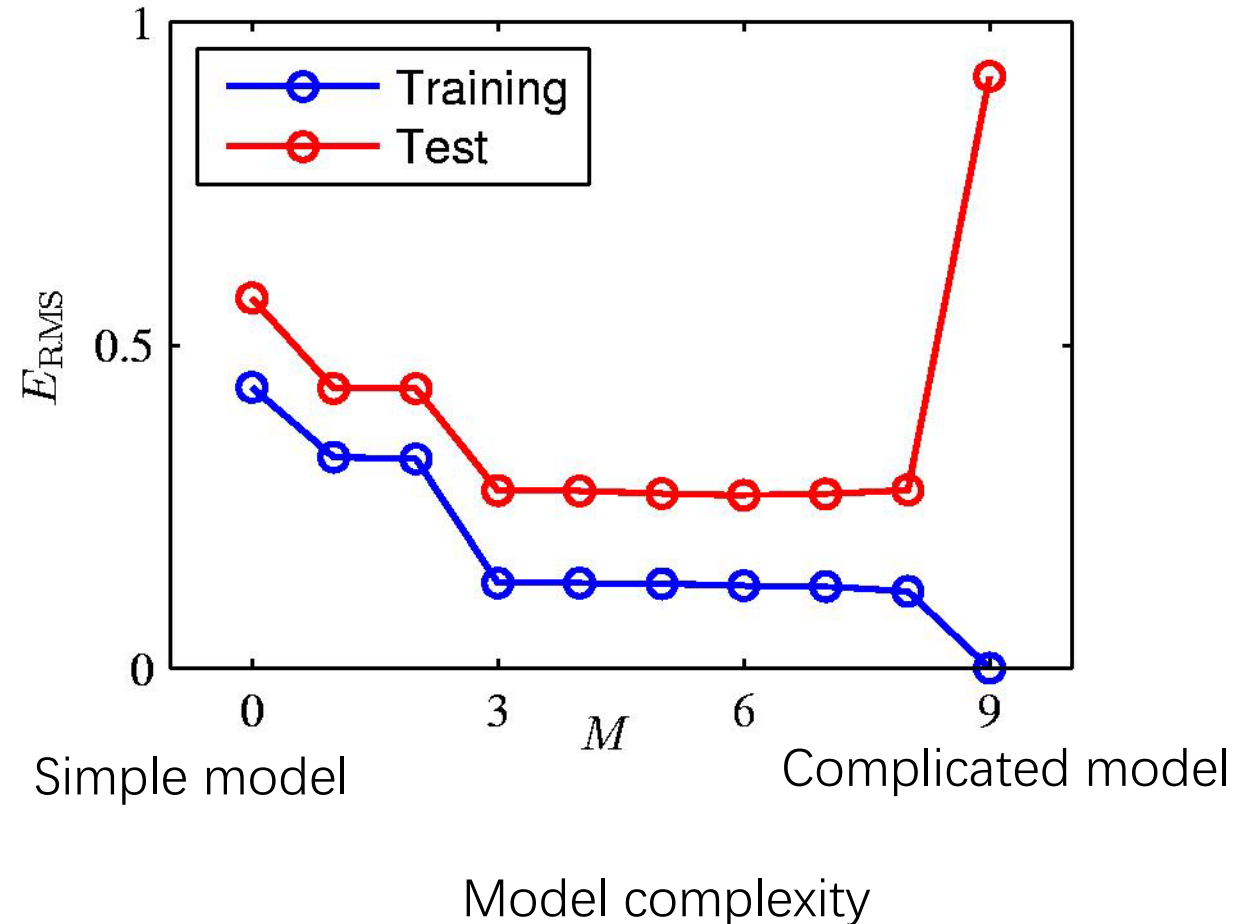
	$M = 0$	$M = 1$	$M = 3$	$M = 9$
w_0^*	0.19	0.82	0.31	0.35
w_1^*		-1.27	7.99	232.37
w_2^*			-25.43	-5321.83
w_3^*			17.37	48568.31
w_4^*				-231639.30
w_5^*				640042.26
w_6^*				-1061800.52
w_7^*				1042400.18
w_8^*				-557682.99
w_9^*				125201.43

Ridge Regression

- How to control the size of the coefficients in Regression?
- $\mathbf{w}^* = \operatorname{argmin} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \mathbf{w})^2 + \lambda \sum_{j=1}^p w_j^2$
- Local smoothness

Regularization

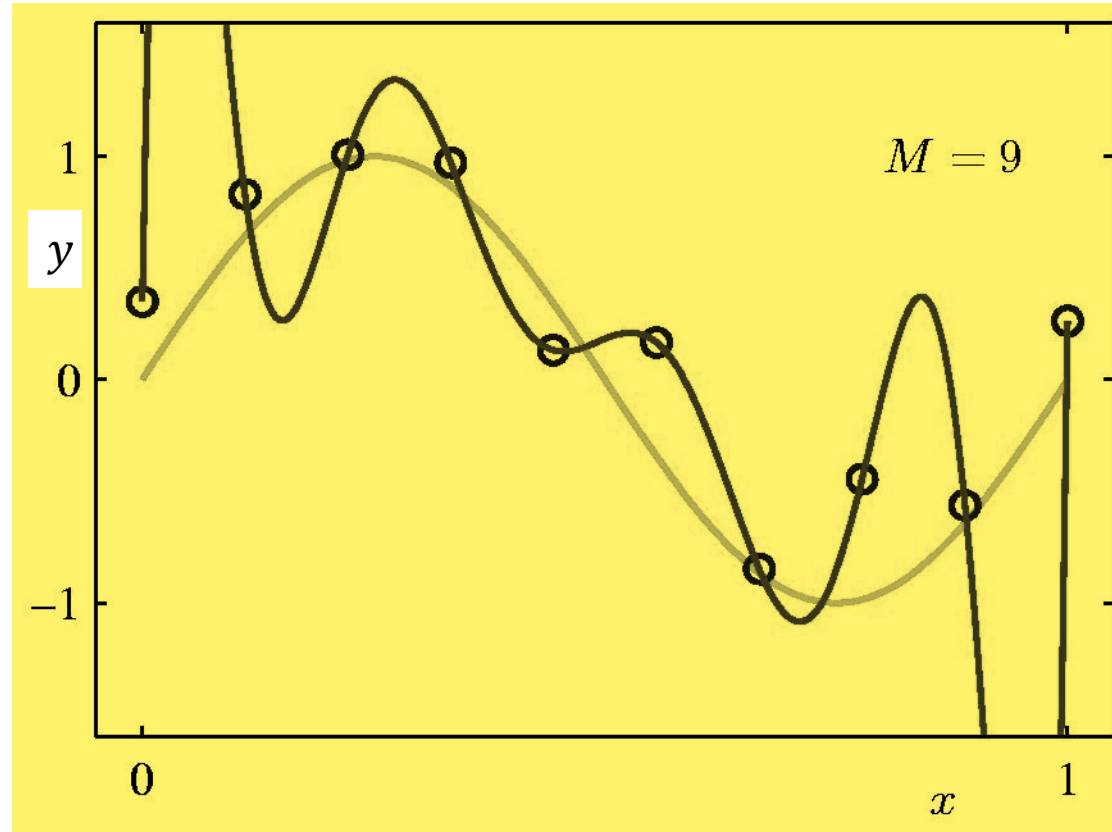
- Some of our priors to the model
- Reduce the size of hypothesis space
- Some good models may not be learned, but avoid learning some bad models



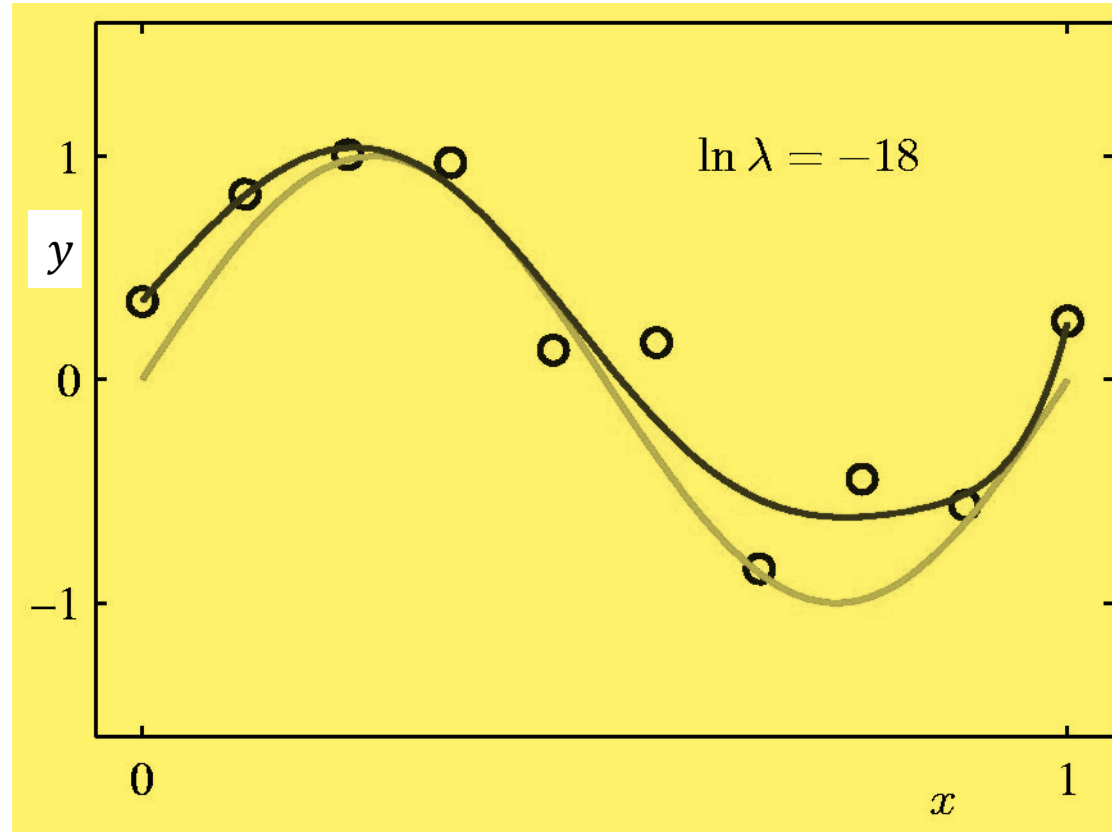
Polynomial Coefficients

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

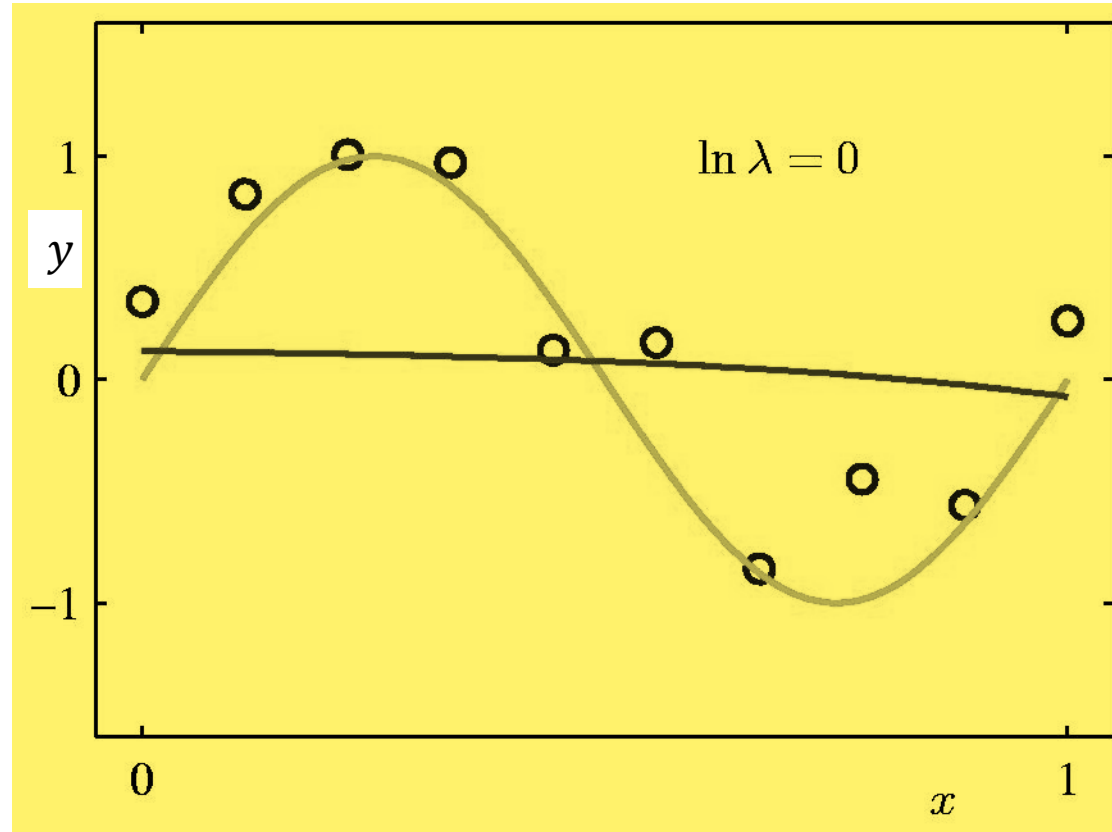
9th Order Polynomial



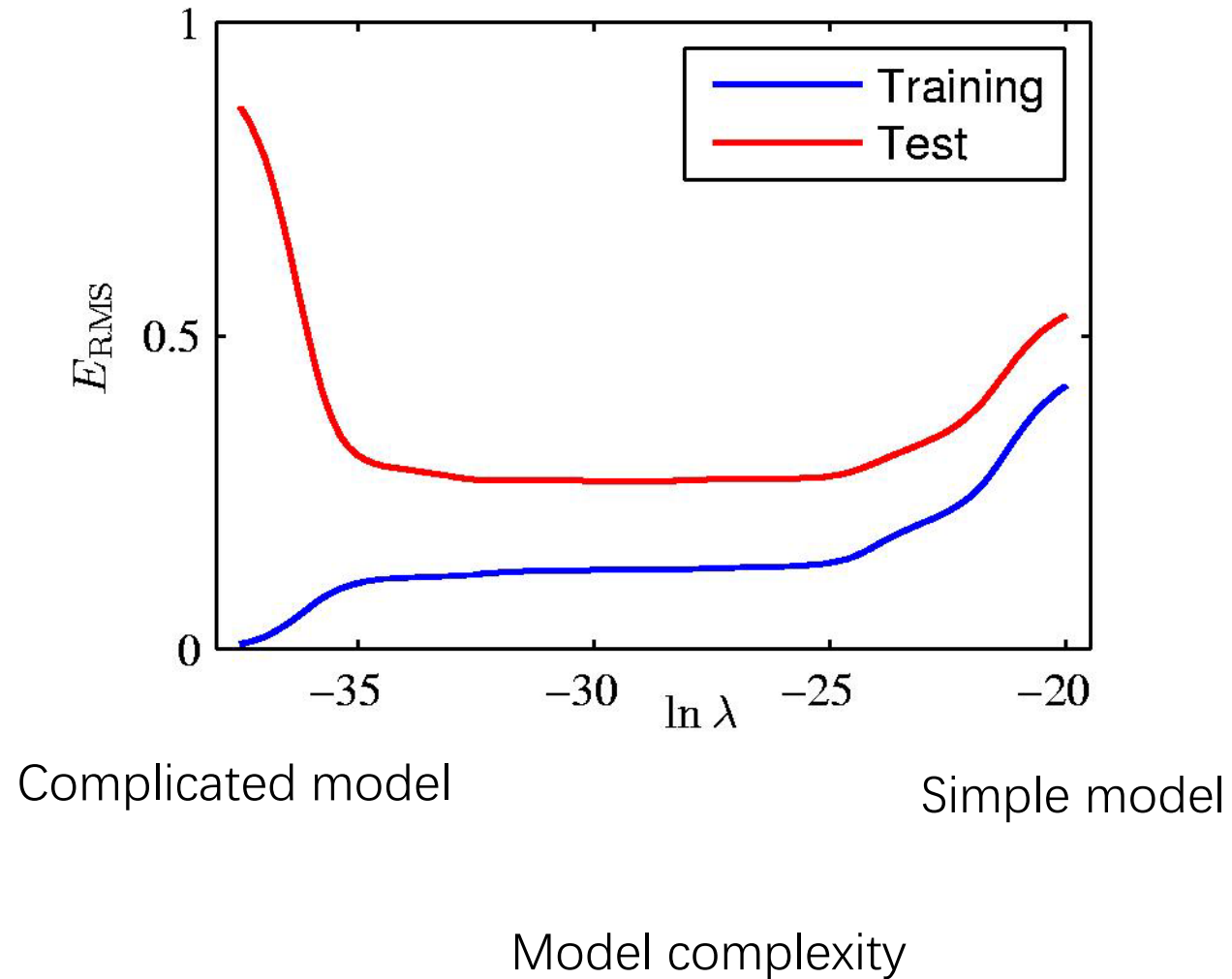
Regularization ($\ln \lambda = -18$)



Regularization ($\ln \lambda = 0$)



Regularization



Bias & Variance Decomposition

The Bias-Variance Decomposition

$$\begin{aligned}(y - f(\mathbf{x}))^2 &= (y - E(y|\mathbf{x}) + E(y|\mathbf{x}) - f(\mathbf{x}))^2 \\ &= (y - E(y|\mathbf{x}))^2 + (E(y|\mathbf{x}) - f(\mathbf{x}))^2 + 2(y - E(y|\mathbf{x}))(E(y|\mathbf{x}) - f(\mathbf{x}))\end{aligned}$$

- Expected Prediction Error:

$$\text{EPE}(f) = \int (f(\mathbf{x}) - E(y|\mathbf{x}))^2 p(\mathbf{x}) d\mathbf{x} + \int \text{var}(y|\mathbf{x}) p(\mathbf{x}) d\mathbf{x}$$

$$[f(\mathbf{x}; D) - E_D(f(\mathbf{x}; D)) + E_D(f(\mathbf{x}; D)) - E(y|\mathbf{x})]^2$$

$$(f(\mathbf{x}; D) - E_D(f(\mathbf{x}; D)))^2 + (E_D(f(\mathbf{x}; D)) - E(y|\mathbf{x}))^2 + 2(f(\mathbf{x}; D) - E_D(f(\mathbf{x}; D)))(E_D(f(\mathbf{x}; D)) - E(y|\mathbf{x}))$$

$$\underbrace{E_D \{ [f(\mathbf{x}; D) - E_D(f(\mathbf{x}; D))]^2 \}}_{\text{Variance}} + \underbrace{\{ E_D(f(\mathbf{x}; D)) - E(y|\mathbf{x}) \}^2}_{(\text{Bias})^2}$$

Variance

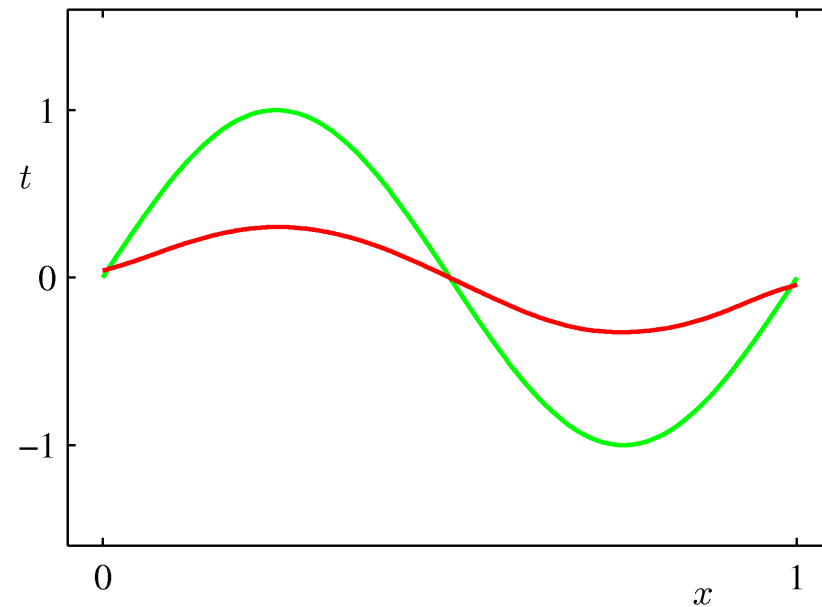
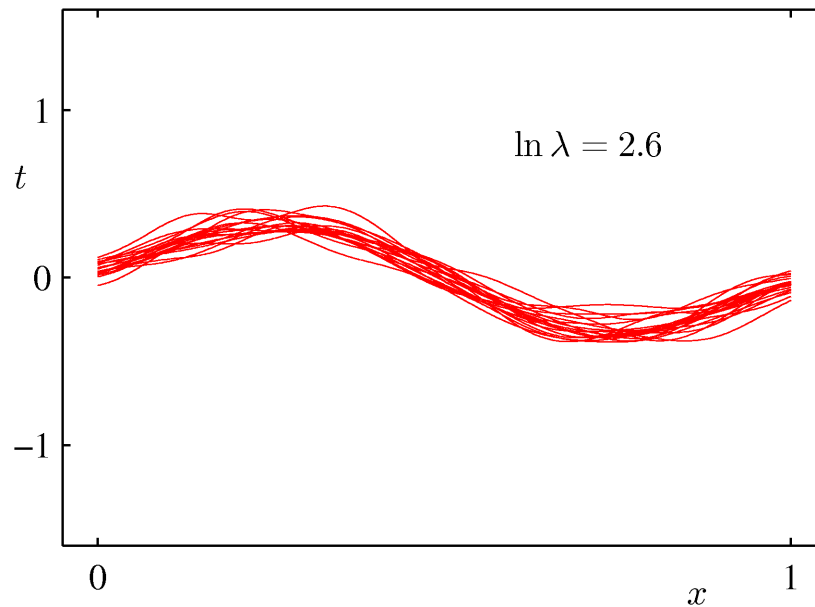
(Bias)²

The Bias-Variance Decomposition

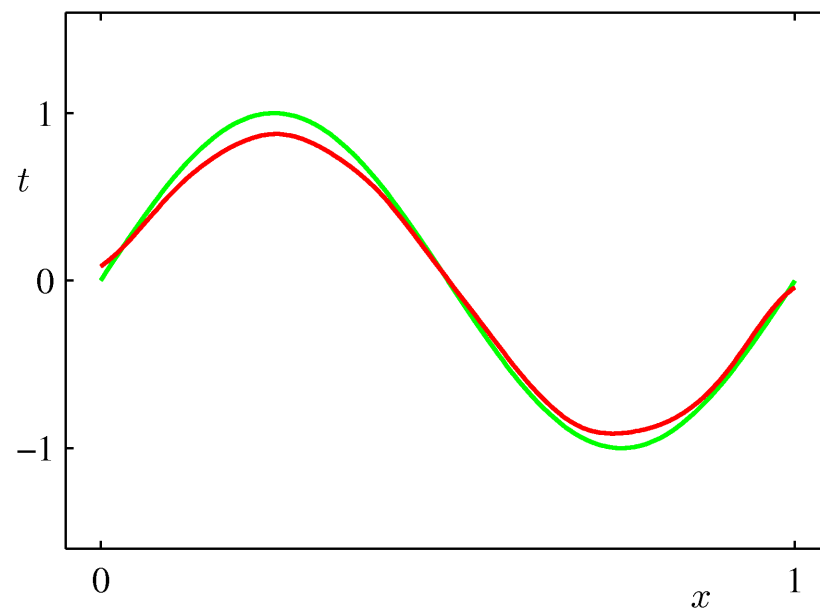
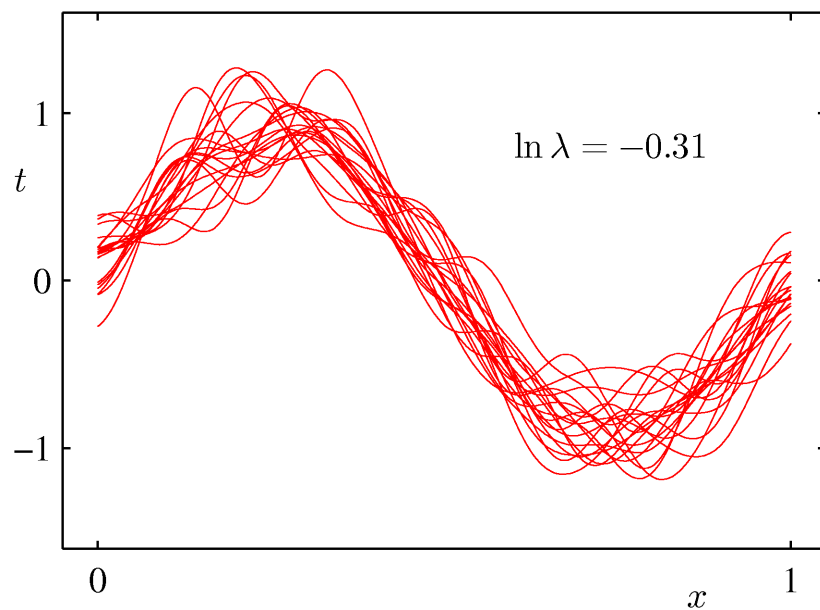
- Through some derivations.....
- Prediction Error = bias + variance + noise

The Bias-Variance Decomposition

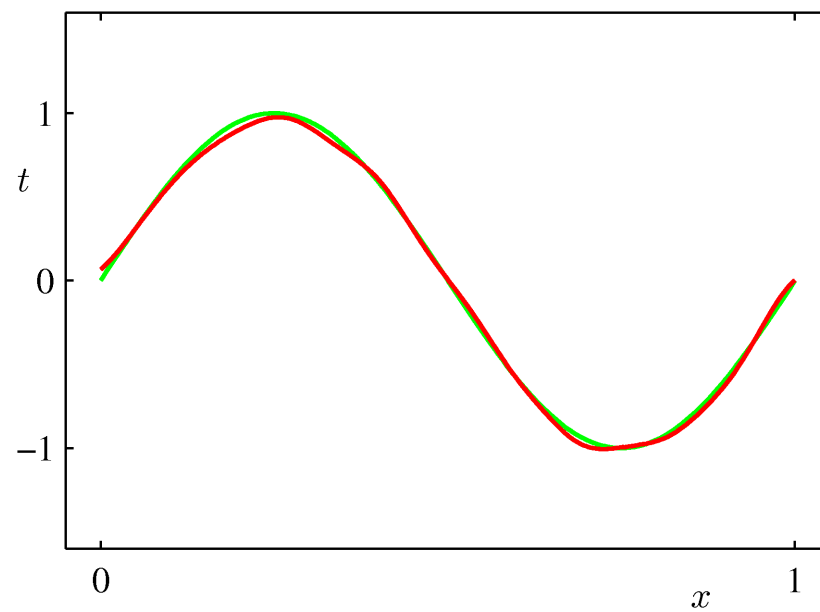
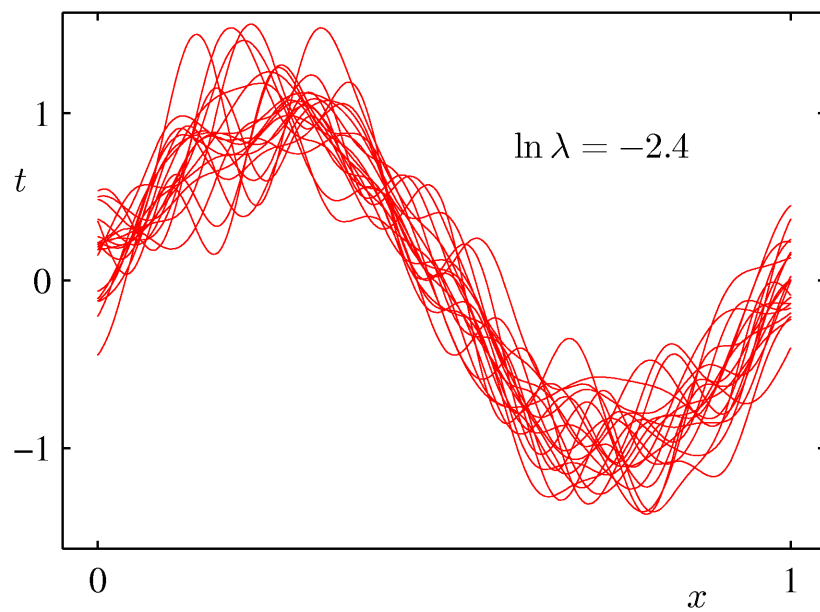
- Example: varying the degree of regularization for some dataset, λ



The Bias-Variance Decomposition

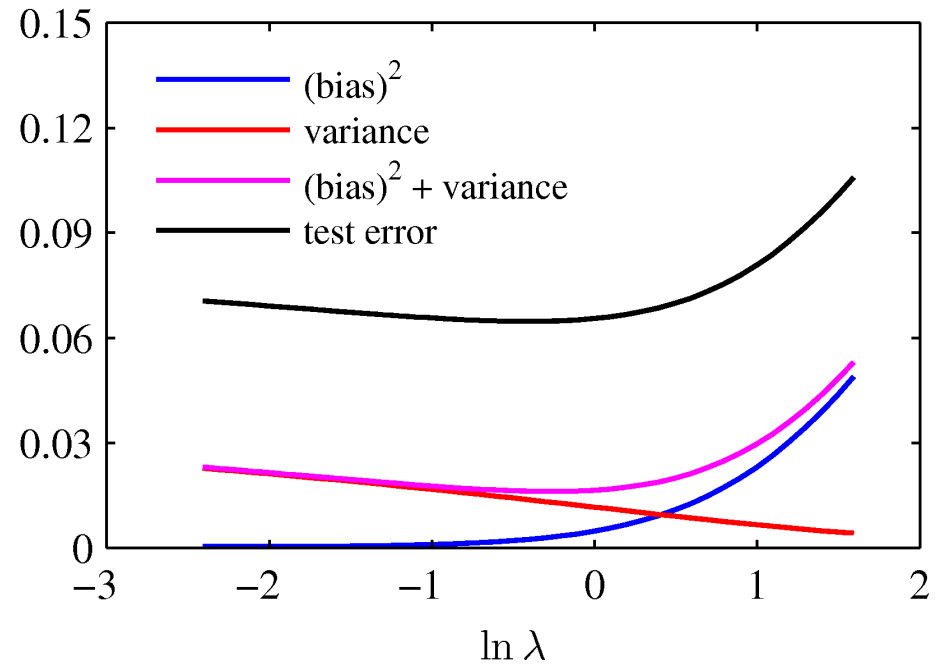


The Bias-Variance Decomposition



The Bias-Variance Trade-off

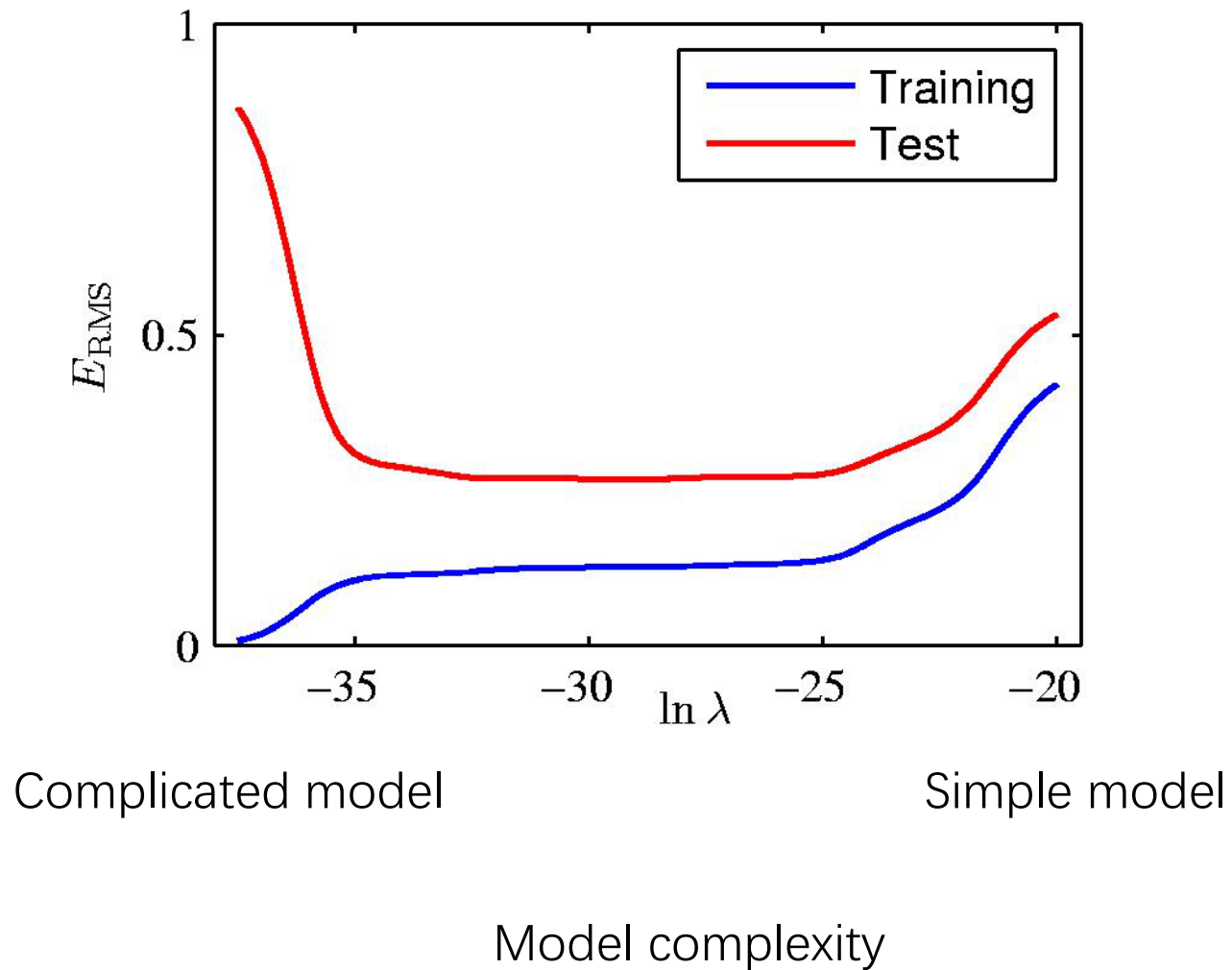
- Over-regularized model (large λ) \rightarrow high bias
- Under-regularized model (small λ) \rightarrow high variance.



Underfitting VS Overfitting

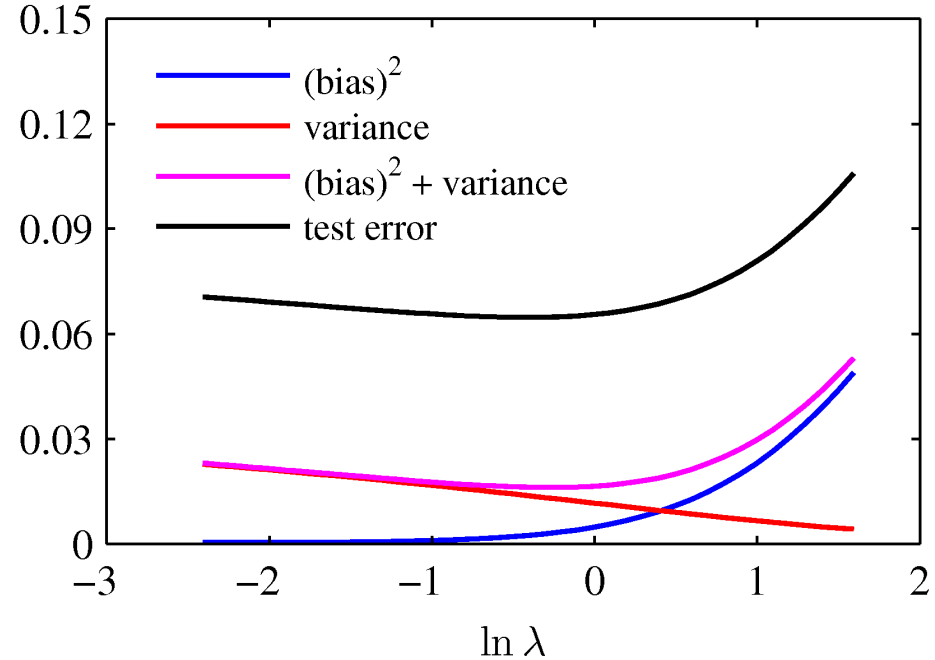
- high bias, low variance → underfitting
- high variance, low bias → overfitting
- low bias, low variance → what we want
- Bias: the training error
- Variance: the difference between training error and test error

Underfitting VS Overfitting



How to avoid underfitting/overfitting?

- Underfitting: lower the bias
 - Add more features
 - Use more complex models
 - Decrease the regularization part
- Overfitting: lower the variance
 - Decrease some feature
 - Use simpler models
 - **Add more regularizations**
 - Add more data

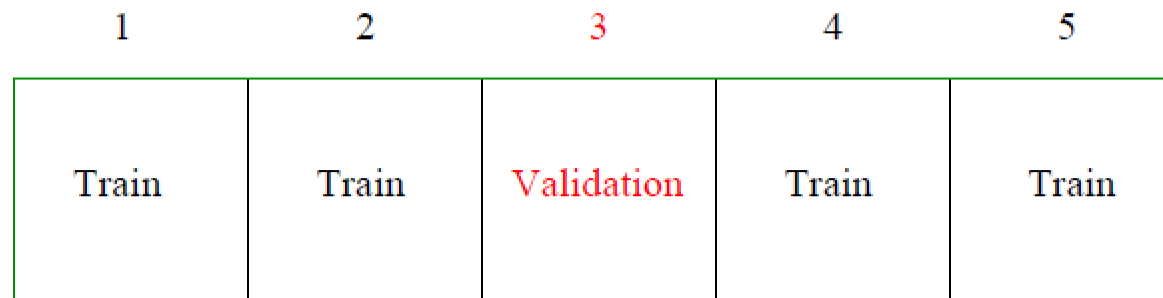


Make choice against the test data

- There are many choices in ML models
 - Features
 - Regularization terms
 - Learning rate
 - Different models
 -
- How to decide which choice is the best?
 - Look at the training error?
 - Look at the testing error?
- If we make a lot of choices against the test data
 - Can the test performance represent the real-world performance?
 - Overfit the test data!

Train Val Test Split

- Split the dataset into three parts:
 - Train: train the model
 - Val: make choices
 - tune hyperparameters
 - choose models
 - choose features
 - Test: evaluate the final model
- K-Fold Cross-Validation



Question?

Thanks and welcome to give us suggestions and feedbacks afterwards.