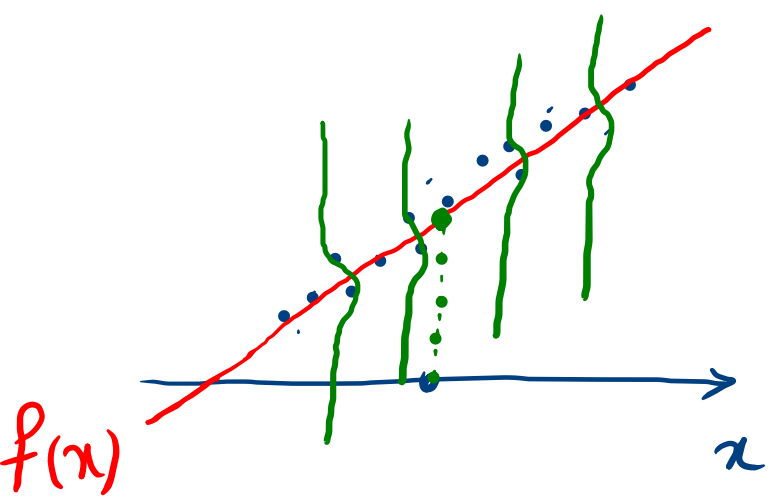


$$g(x) = \underline{\underline{\omega}}^T x$$

?

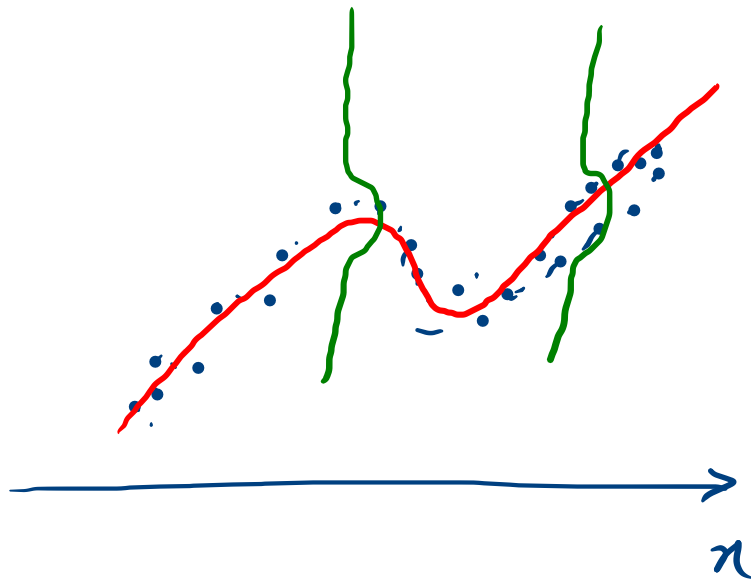
$$\omega^* = \arg \min_{\omega} \underbrace{\sum_{i=1}^n (\omega^T x_i - y_i)^2}_{\text{SSE}}$$

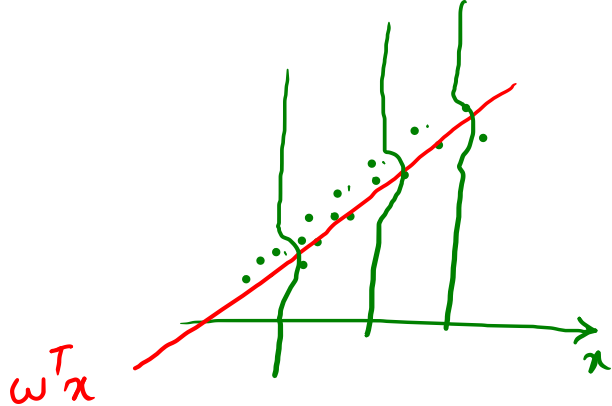
Least Squares



$$y = \underline{f(x)} + \textcircled{n}$$

$$\boxed{n \sim \mathcal{N}(0, \textcircled{\sigma^2})}$$





$$\begin{cases} y = \underline{w^T x} + \epsilon \\ \epsilon \sim \mathcal{N}(0, \sigma^2) \end{cases}$$

$$P(y|x) = \mathcal{N}(y; \underline{w^T x}, \sigma^2) \leftarrow$$

$$D = \{(\overset{\downarrow}{x_1}, y_1), \dots, (x_n, y_n)\}$$

$$E[y] = E[y|x] = E[\underline{w^T x} + \epsilon | x]$$

$$= w^T x$$

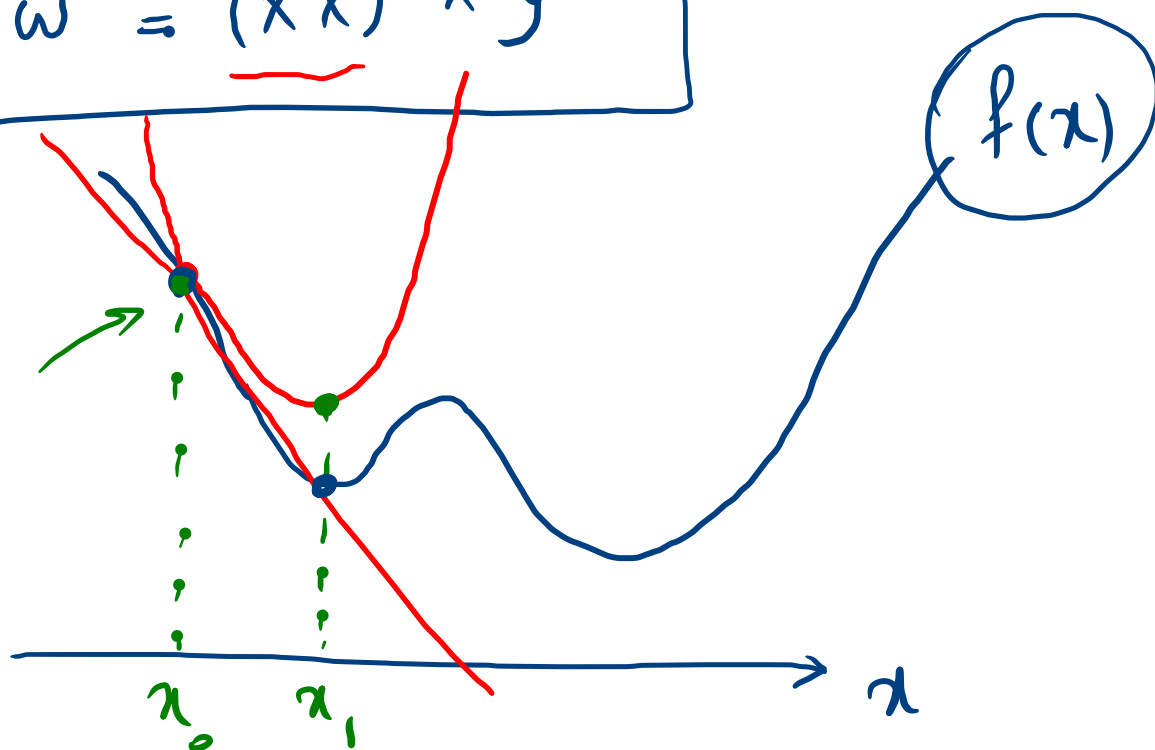
$$w_{ML}^* = \arg \max_w \ln P(D|w) = \arg \max_w \ln P(y_1, \dots, y_n | x_1, \dots, x_n, w) \quad \text{L}(w)$$

$$\stackrel{\text{i.i.d.}}{=} \arg \max_w \sum_{i=1}^n \ln P(y_i | x_i, w) = \arg \max_w \sum_{i=1}^n \ln \mathcal{N}(y_i; w^T x_i, \sigma^2)$$

$$= \arg \max_w \sum_{i=1}^n \ln \left( \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{1}{2} \frac{(w^T x_i - y_i)^2}{\sigma^2}\right) \right)$$

$$= \arg \max_w \sum_i \underbrace{-\ln \sqrt{2\pi} \sigma}_{\text{SSE } f(w)} - \frac{1}{2} \frac{(w^T x_i - y_i)^2}{\sigma^2} = \arg \min_w \sum_{i=1}^n (w^T x_i - y_i)^2$$

$$\omega^* = \underline{(X^T X)^{-1}} X^T y$$



$$X_{n \times d}$$

$$(X^T X)_{d \times d}$$

$$O(d^3) = O(10^{12})$$

$$\downarrow \quad \downarrow$$

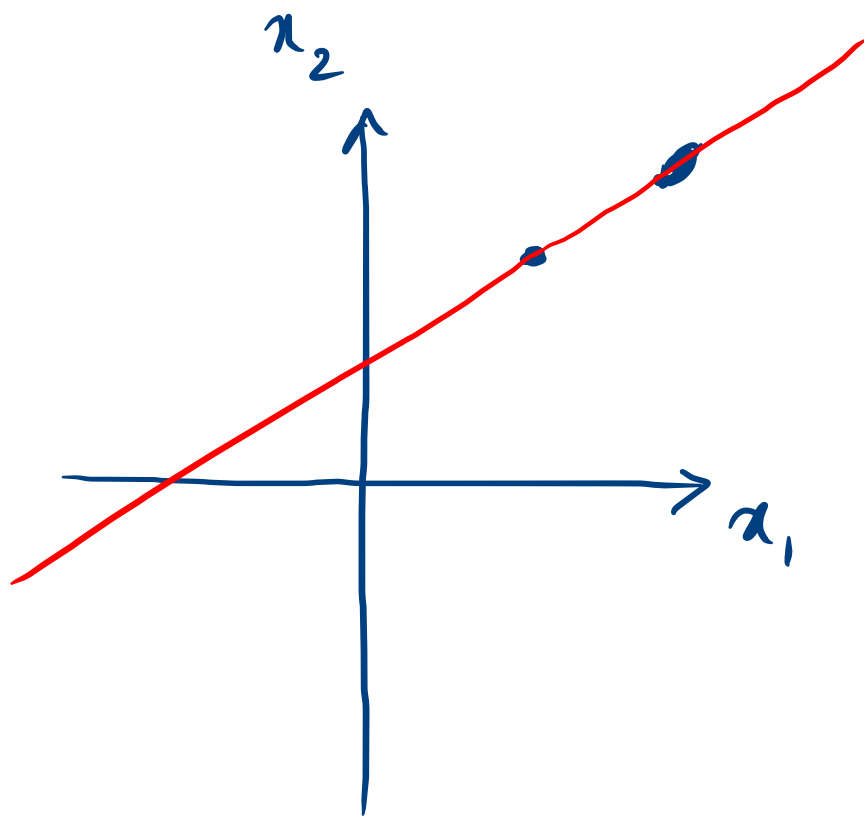
$$ax + b$$

$$a = f'(x_0)$$

$$ax^2 + bx + c$$

$$2ax + b = f'(x_0)$$

$$2a = f''(x_0)$$



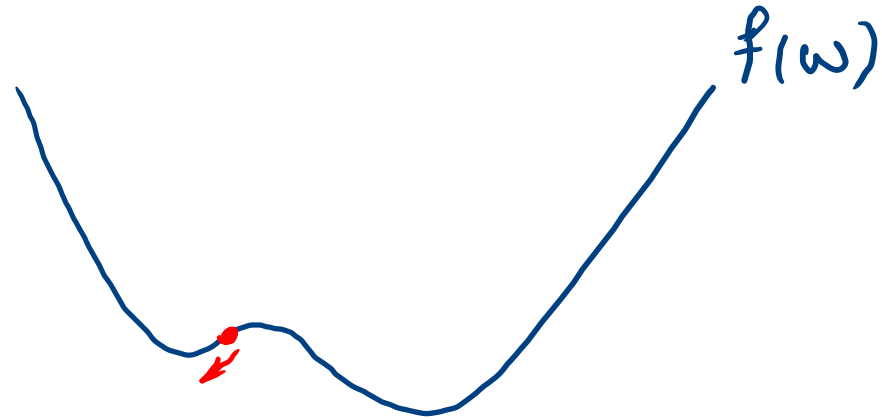
$$d=2$$

$$n=1$$

$$\max_w \log L(w)$$

$$\min_w \text{Loss}(w)$$

Batch



→ stochastic Gradient Descent (SGD)

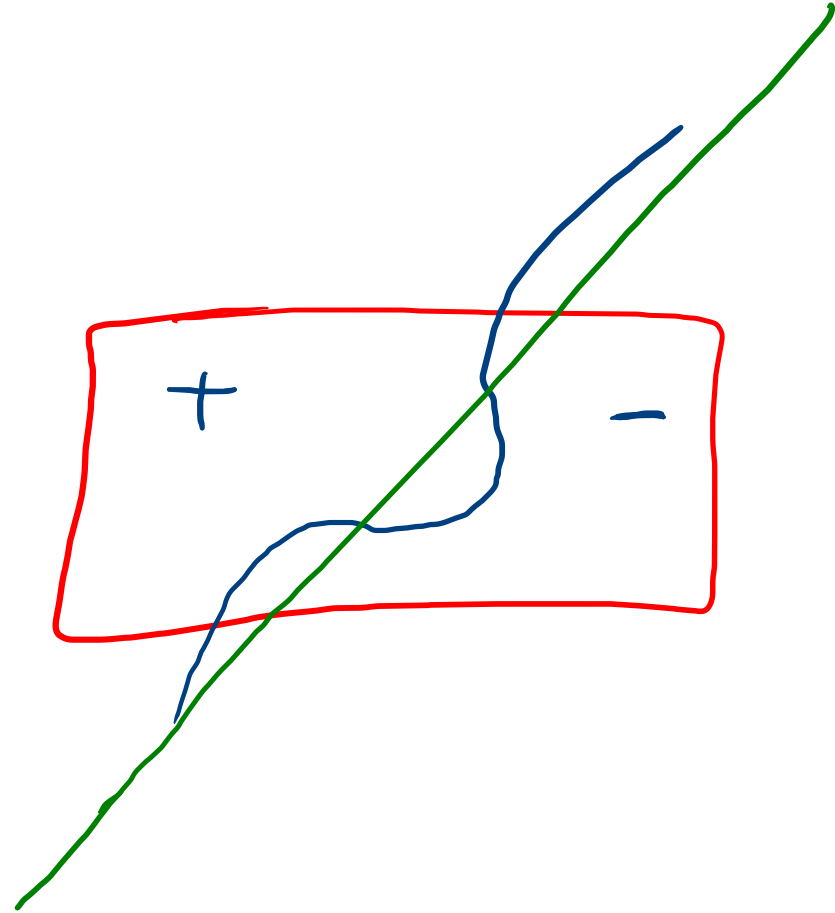
Mini-batch

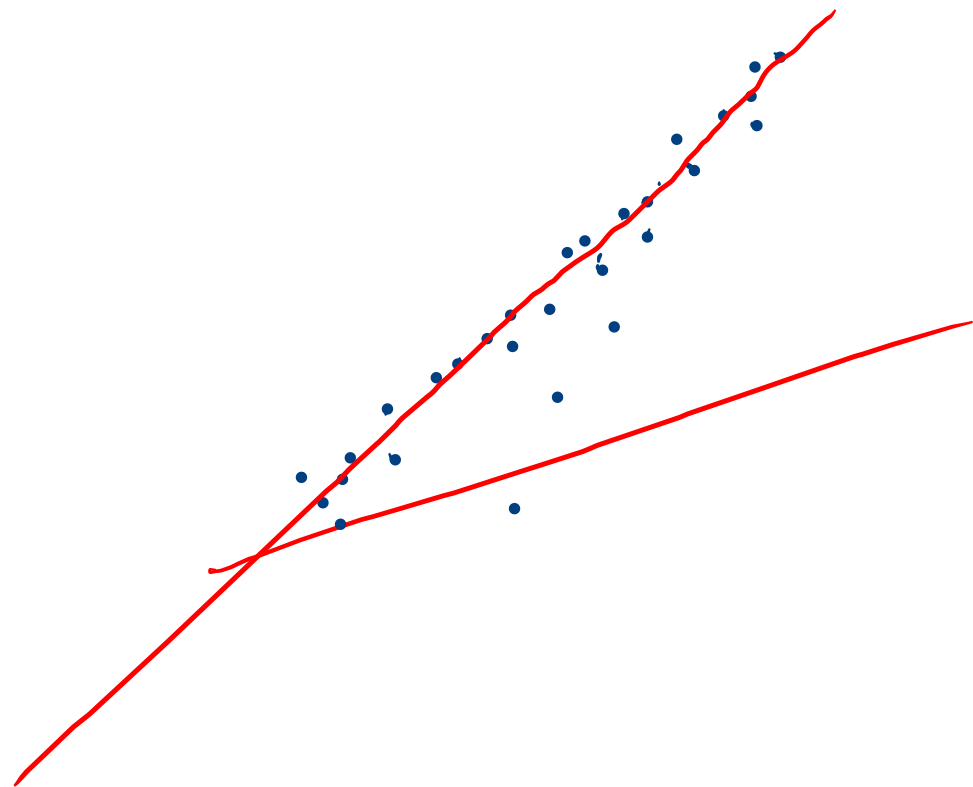
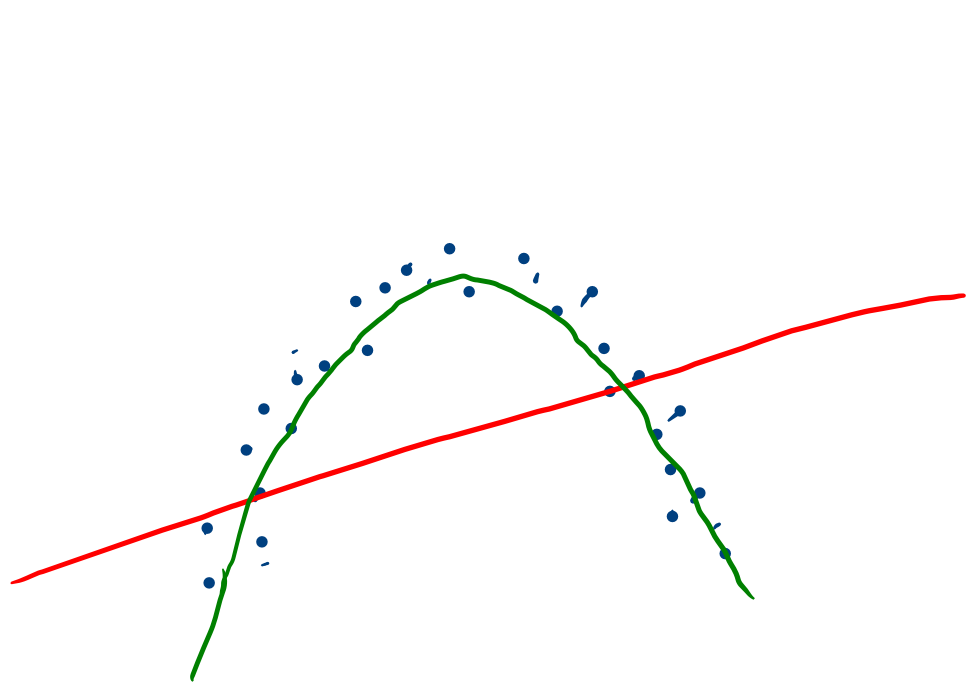
میسر

$$\text{Error} = \text{Bias} + \text{Variance}$$

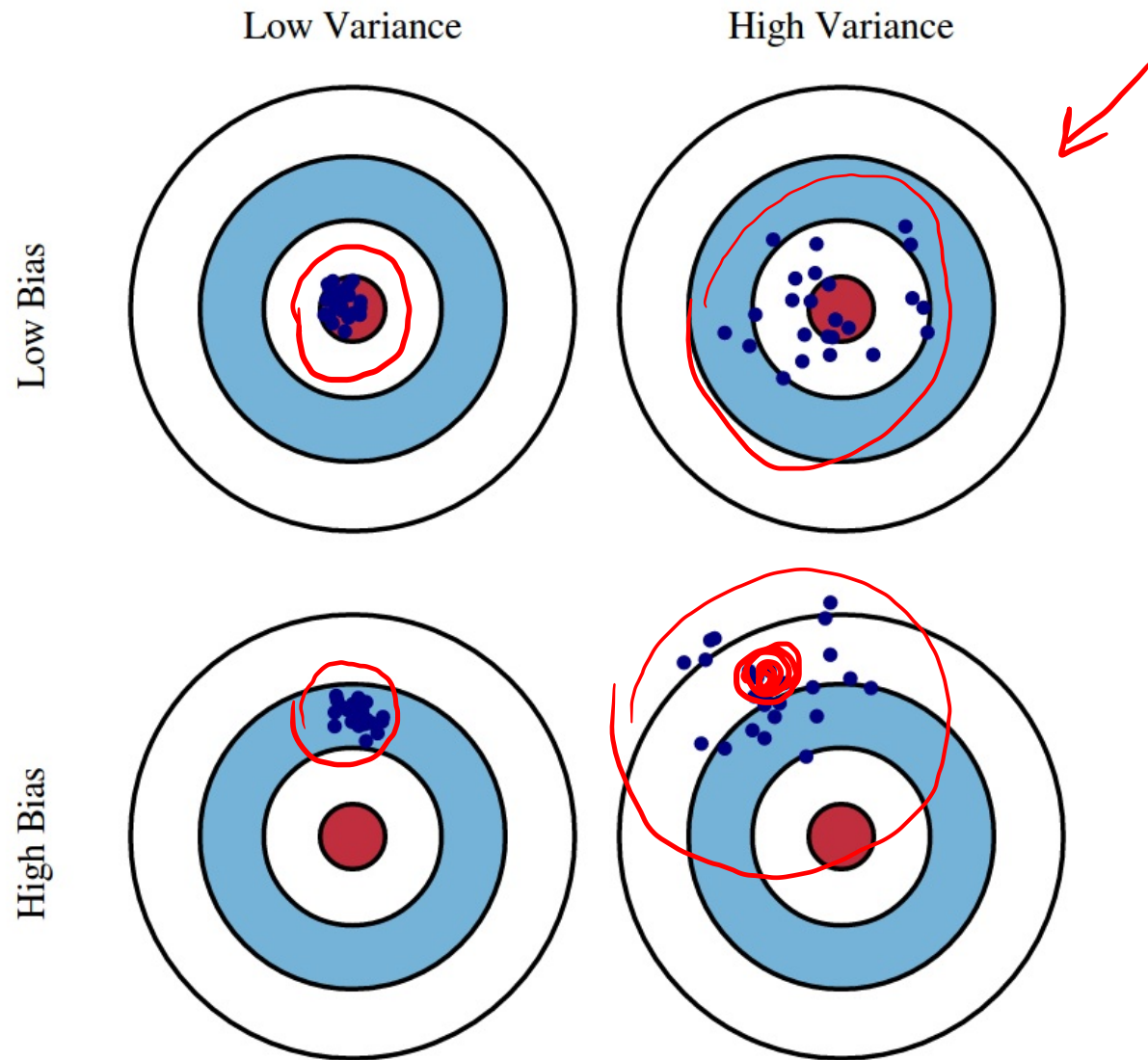
Bias - Variance trade off

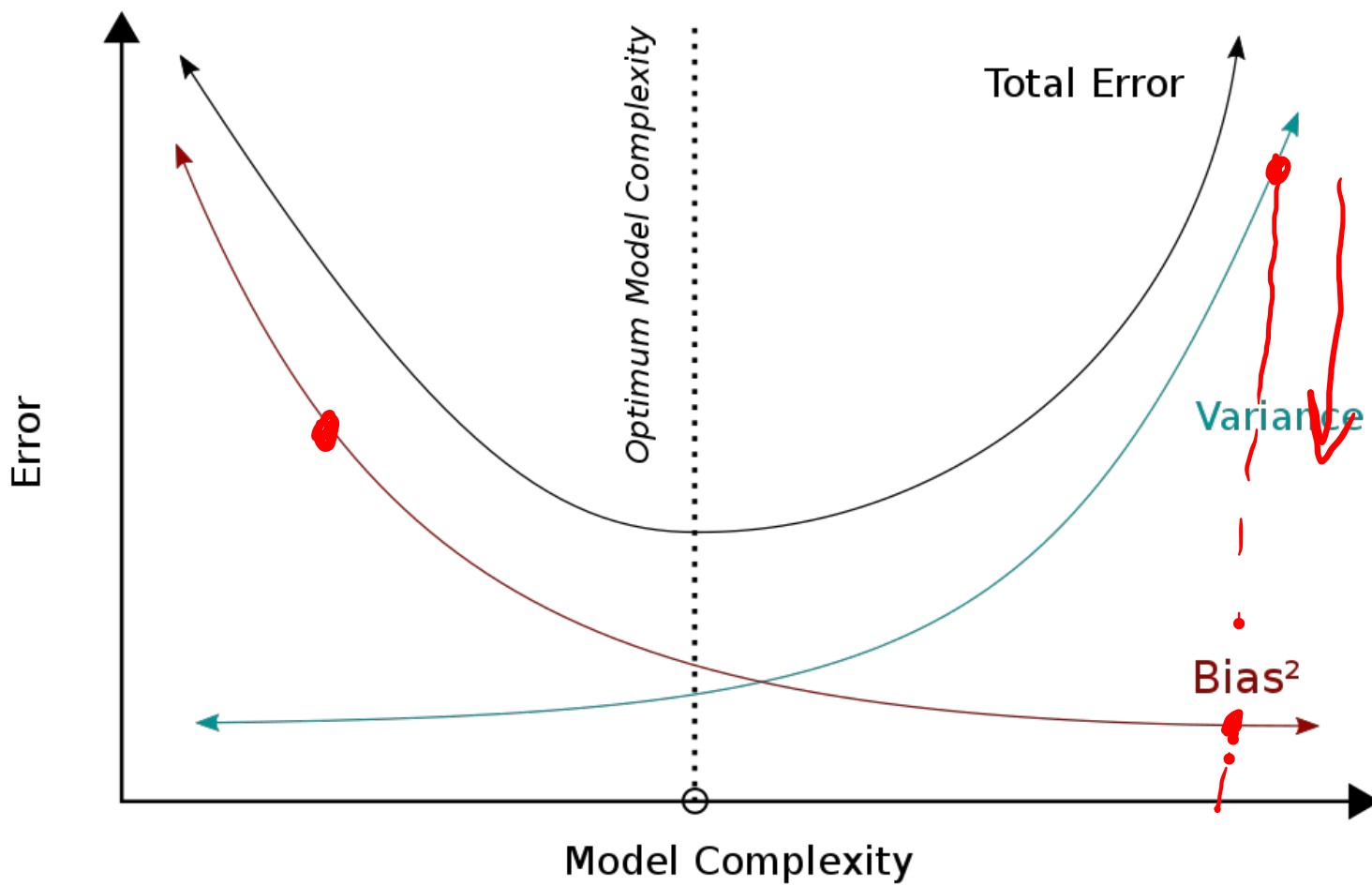
$n \uparrow \Rightarrow \text{variance} \downarrow$

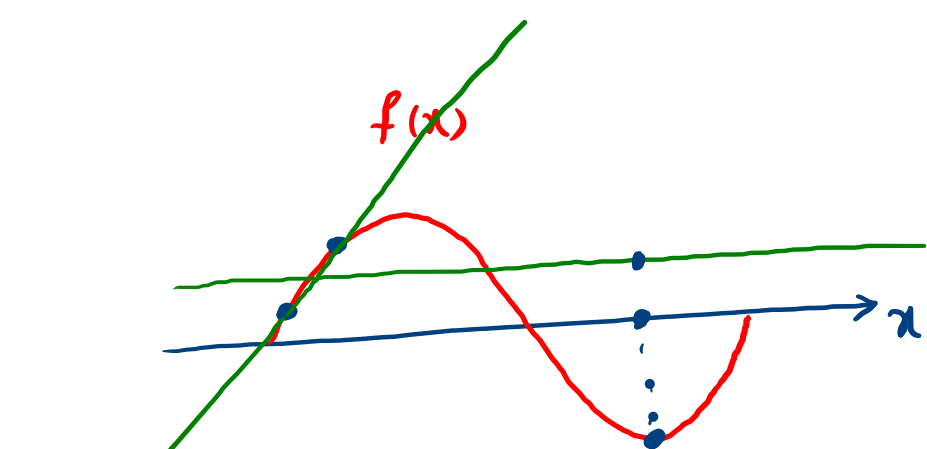




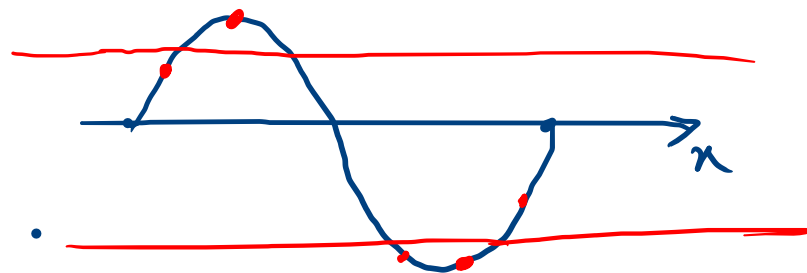
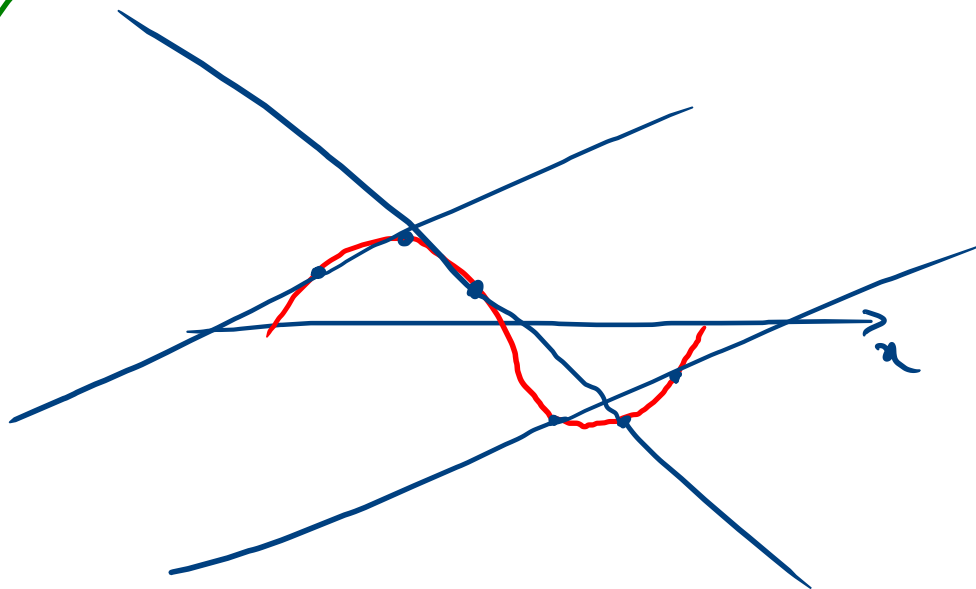






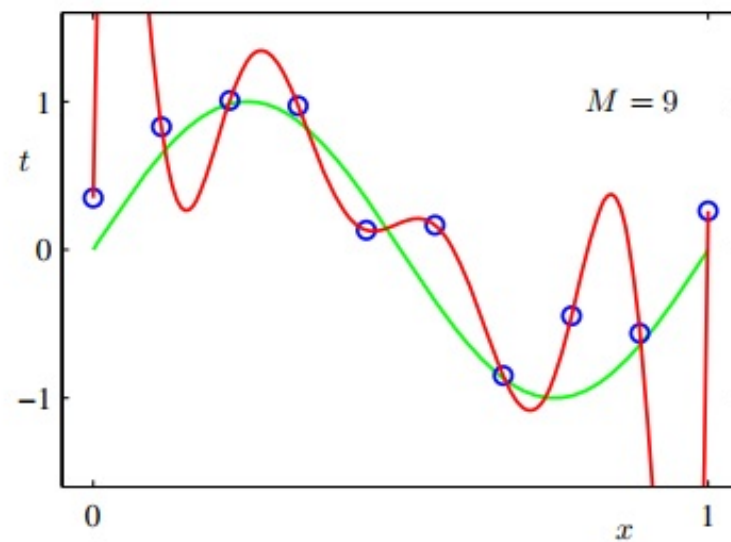
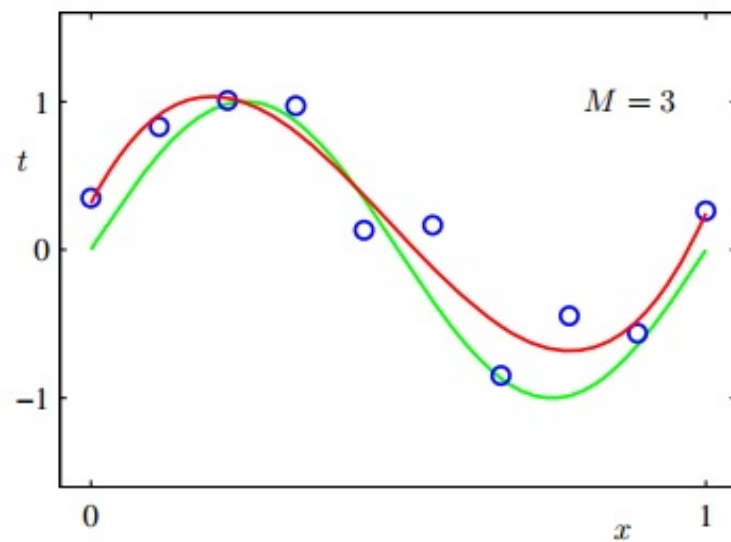
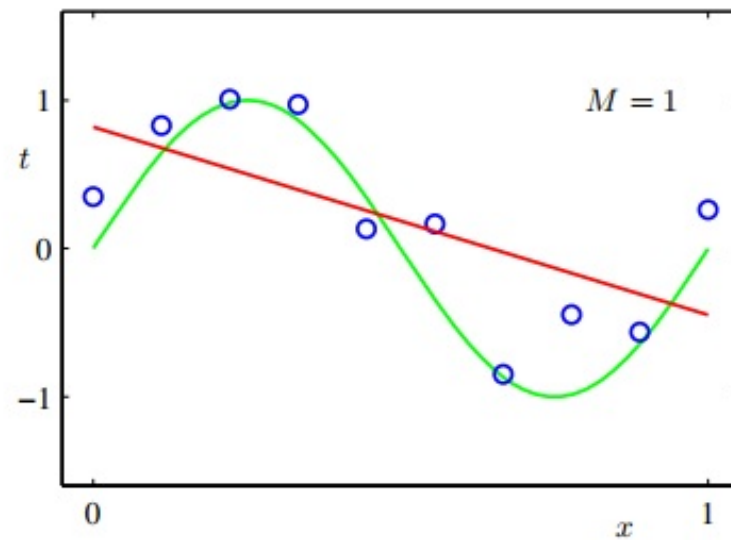
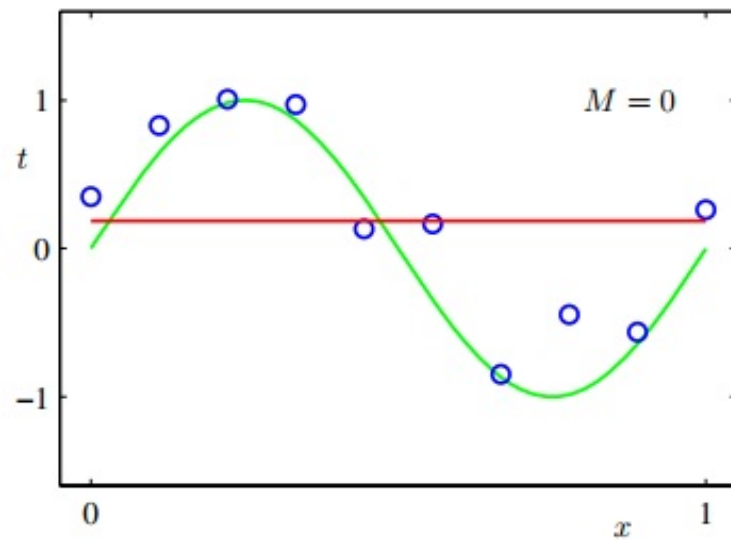


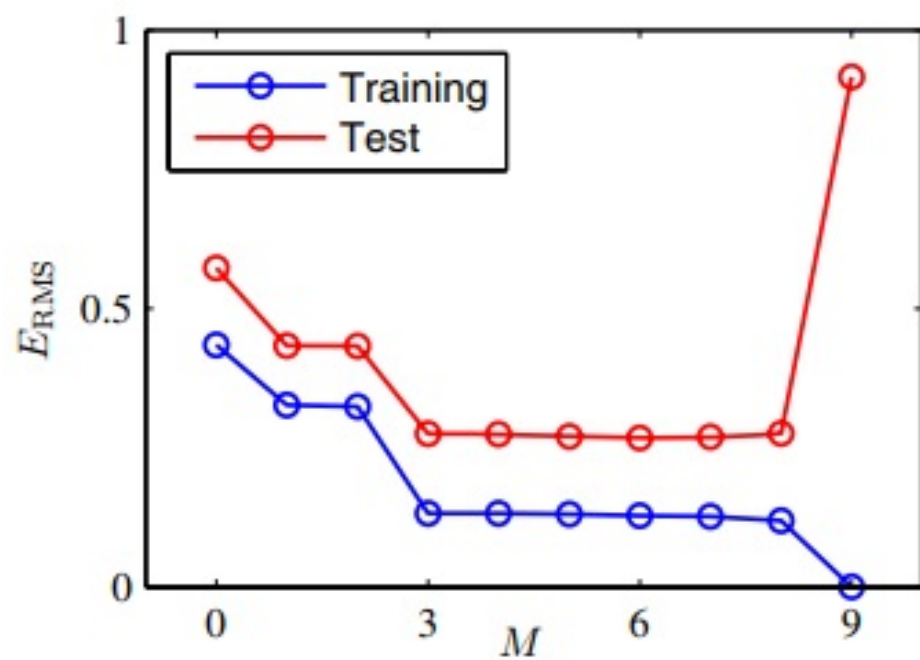
$\rightarrow g(x) = b$  ←  
 $\rightarrow g(x) = ax + b$  ←



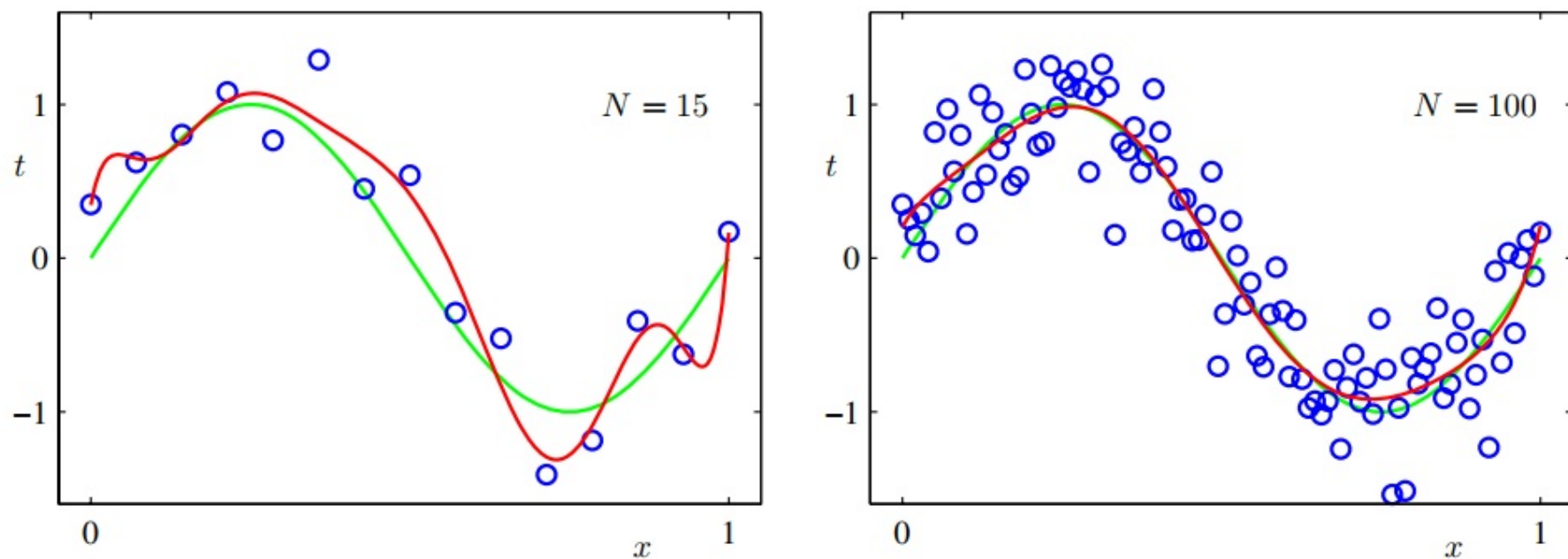
1

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

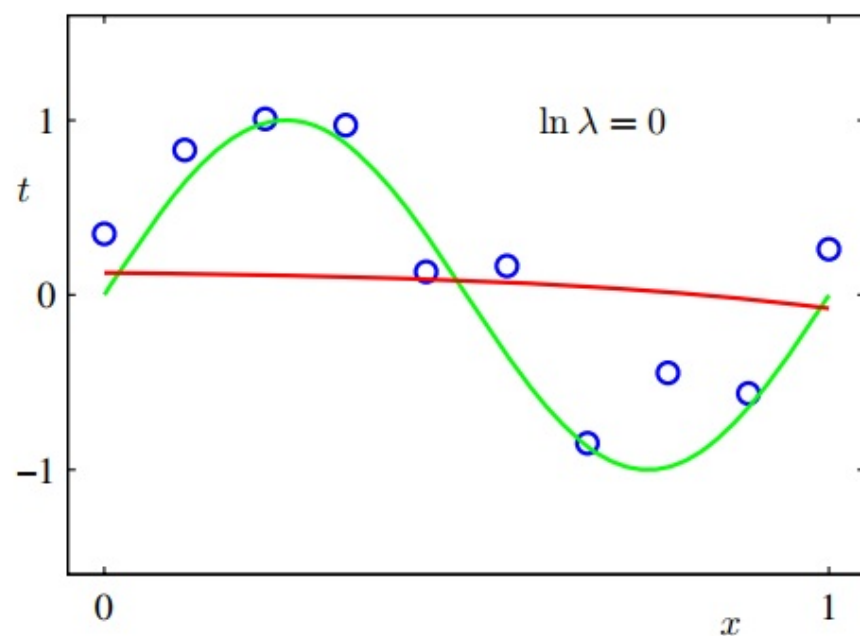
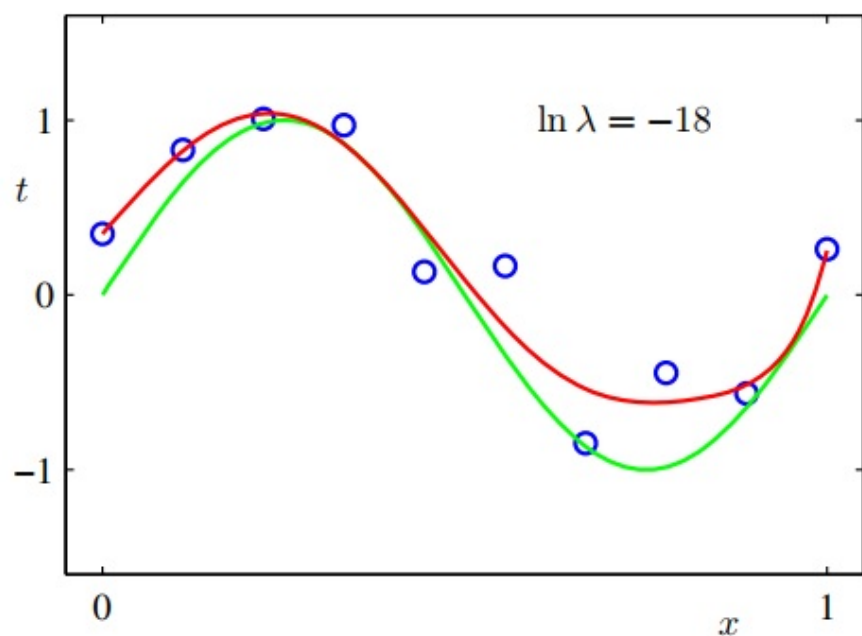




	$M = 0$	$M = 1$	$M = 6$	$M = 9$
$w_0^*$	0.19	0.82	0.31	0.35
$w_1^*$		-1.27	7.99	232.37
$w_2^*$			-25.43	-5321.83
$w_3^*$			17.37	48568.31
$w_4^*$				-231639.30
$w_5^*$				640042.26
$w_6^*$				-1061800.52
$w_7^*$				1042400.18
$w_8^*$				-557682.99
$w_9^*$				125201.43



**Figure 1.6** Plots of the solutions obtained by minimizing the sum-of-squares error function using the  $M = 9$  polynomial for  $N = 15$  data points (left plot) and  $N = 100$  data points (right plot). We see that increasing the size of the data set reduces the over-fitting problem.



**Figure 1.7** Plots of  $M = 9$  polynomials fitted to the data set shown in Figure 1.2 using the regularized error function (1.4) for two values of the regularization parameter  $\lambda$  corresponding to  $\ln \lambda = -18$  and  $\ln \lambda = 0$ . The case of no regularizer, i.e.,  $\lambda = 0$ , corresponding to  $\ln \lambda = -\infty$ , is shown at the bottom right of Figure 1.4.



**Table 1.2** Table of the coefficients  $w^*$  for  $M = 9$  polynomials with various values for the regularization parameter  $\lambda$ . Note that  $\ln \lambda = -\infty$  corresponds to a model with no regularization, i.e., to the graph at the bottom right in Figure 1.4. We see that, as the value of  $\lambda$  increases, the typical magnitude of the coefficients gets smaller.

	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
$w_0^*$	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



**Figure 1.8**

Graph of the root-mean-square error (1.3) versus  $\ln \lambda$  for the  $M = 9$  polynomial.

