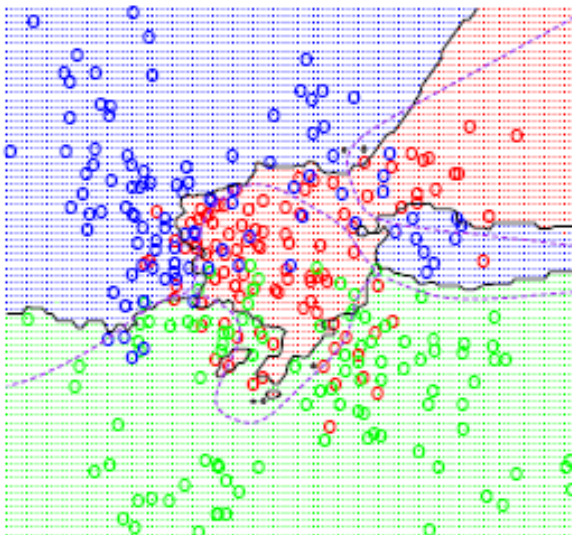


Machine Learning

Lecture 4: Regularization; Bias-Variance tradeoff

The lectures are mainly offered on white board accompanied by some slides.

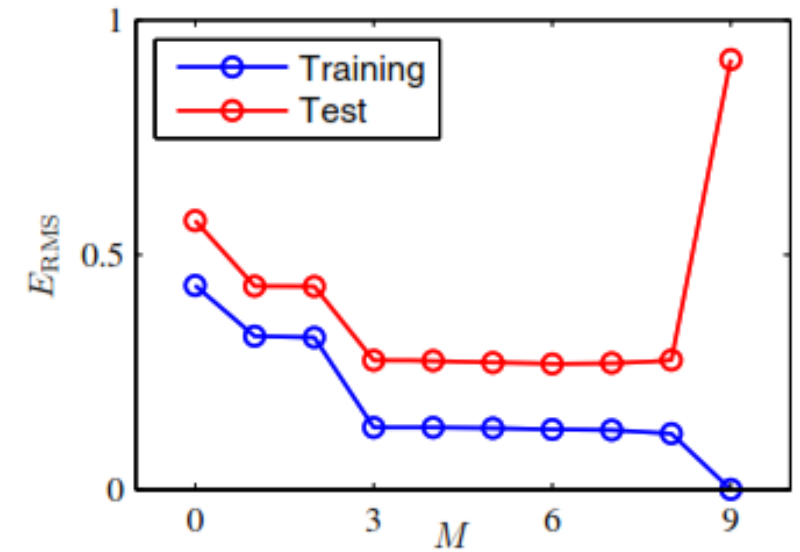
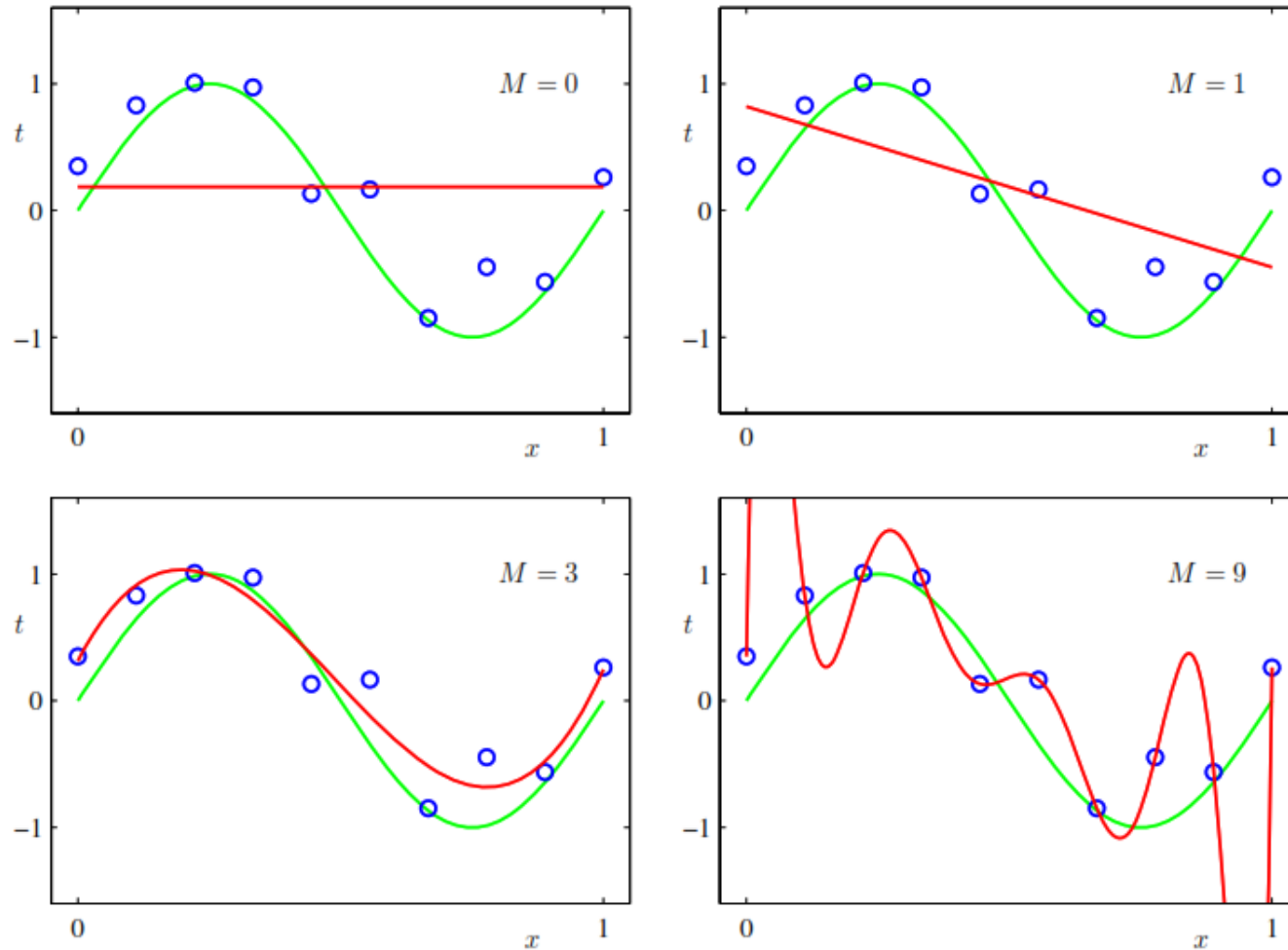


Hesam Montazeri

Department of Bioinformatics, IBB, University of Tehran

Bahman 21, 1398

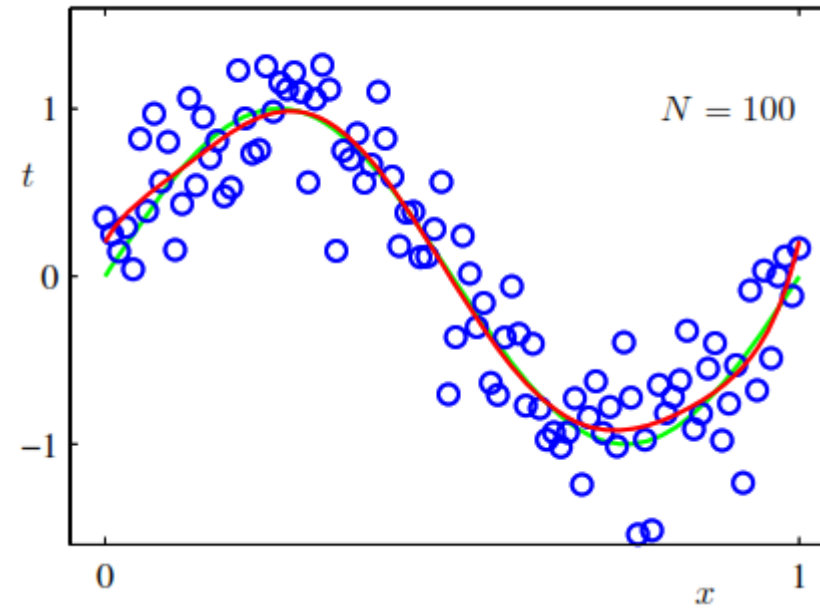
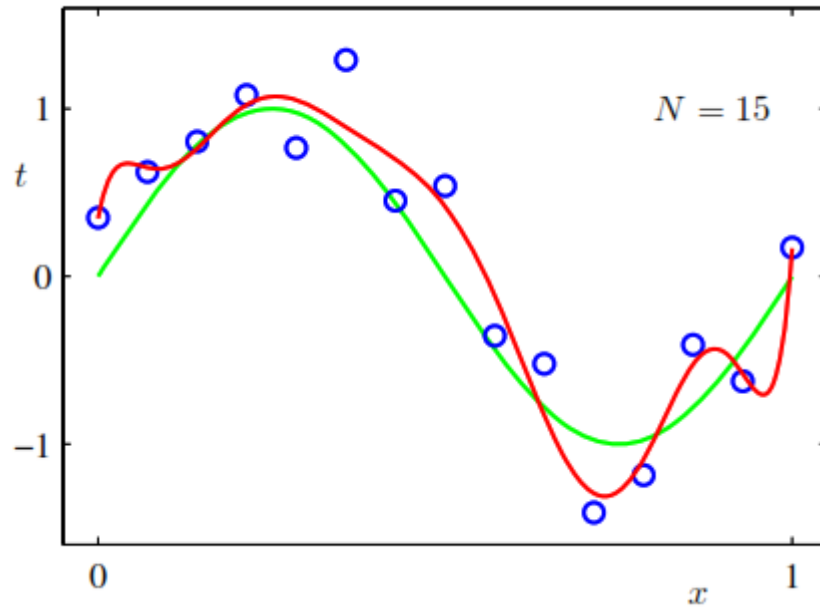
Polynomials having various orders M



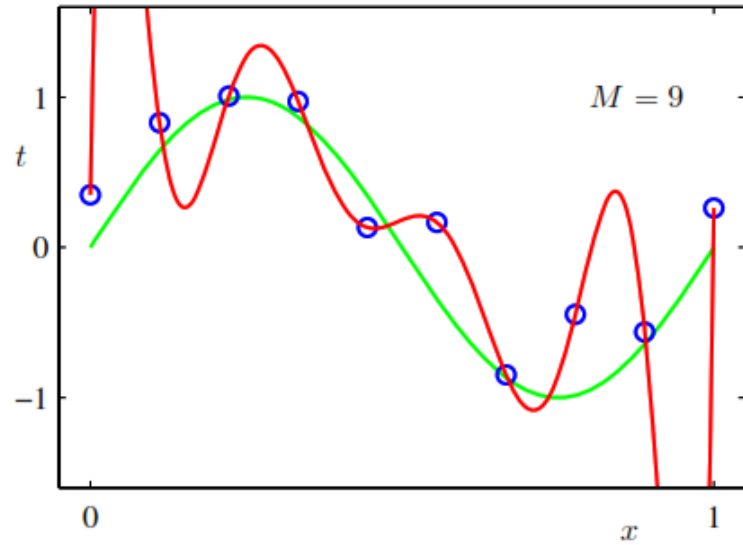
Magnitude of the coefficients increases with p

	$M = 0$	$M = 1$	$M = 6$	$M = 9$
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

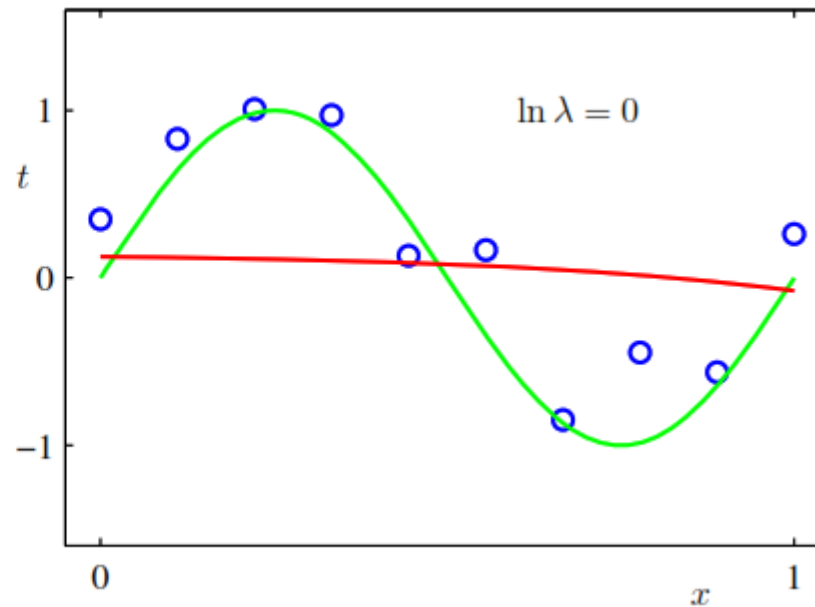
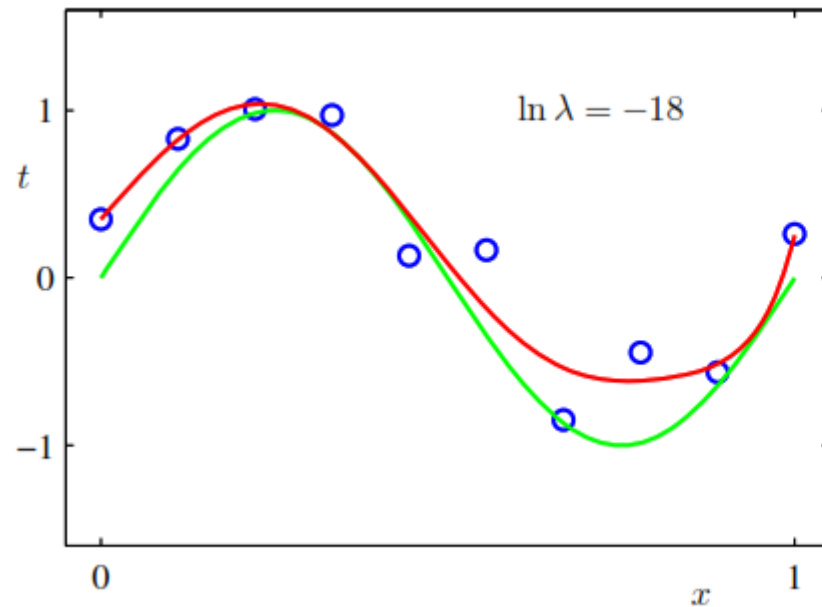
The increasing size of the data set reduces the over-fitting problem

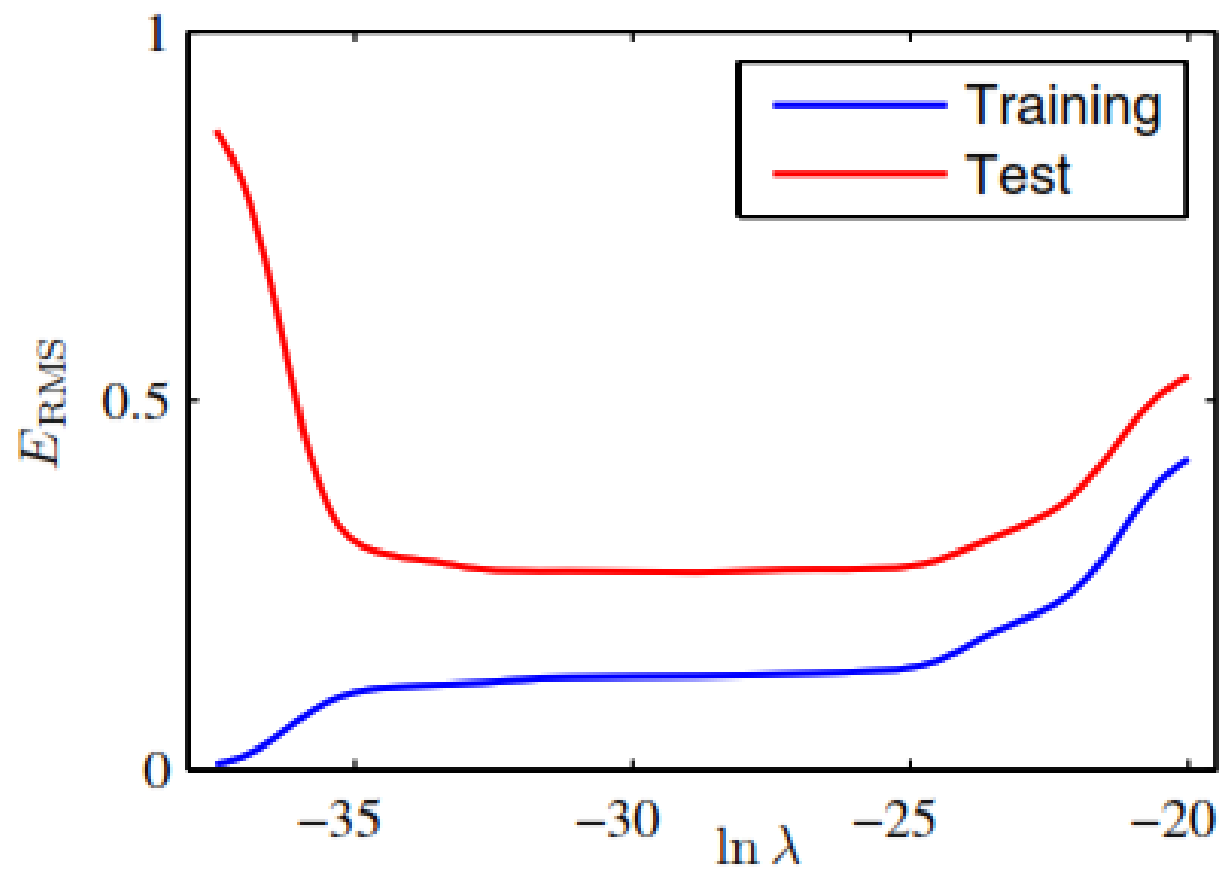


Regularized error function



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





Example: sine target

f

$$f : [-1, 1] \rightarrow \mathbb{R} \quad f(x) = \sin(\pi x)$$

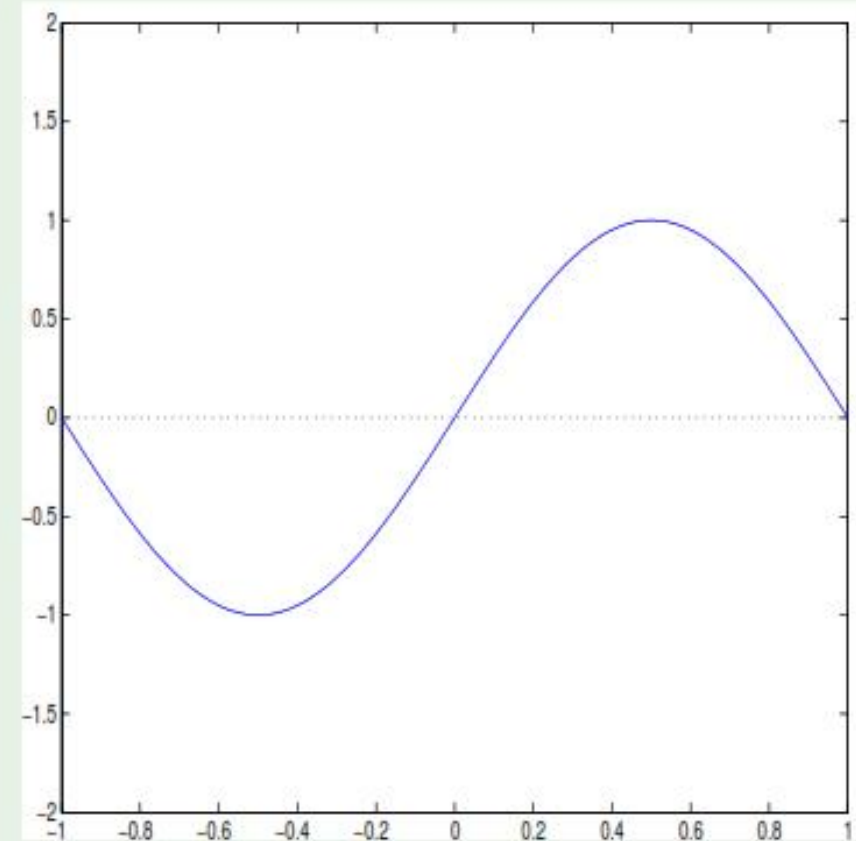
Only two training examples! $N = 2$

Two models used for learning:

$$\mathcal{H}_0: \quad h(x) = b$$

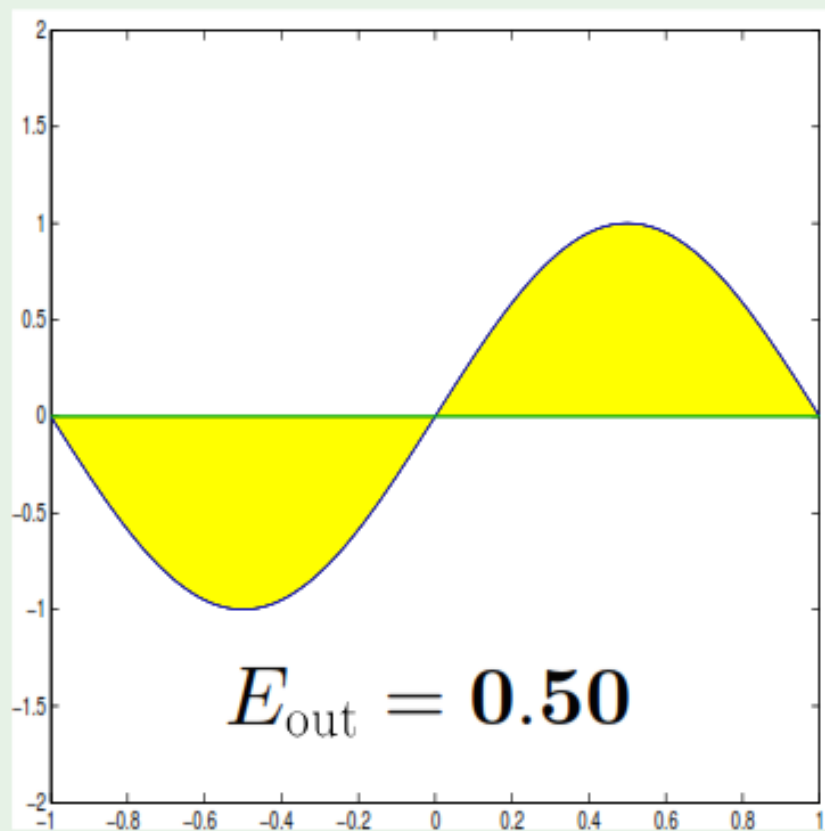
$$\mathcal{H}_1: \quad h(x) = ax + b$$

Which is better, \mathcal{H}_0 or \mathcal{H}_1 ?

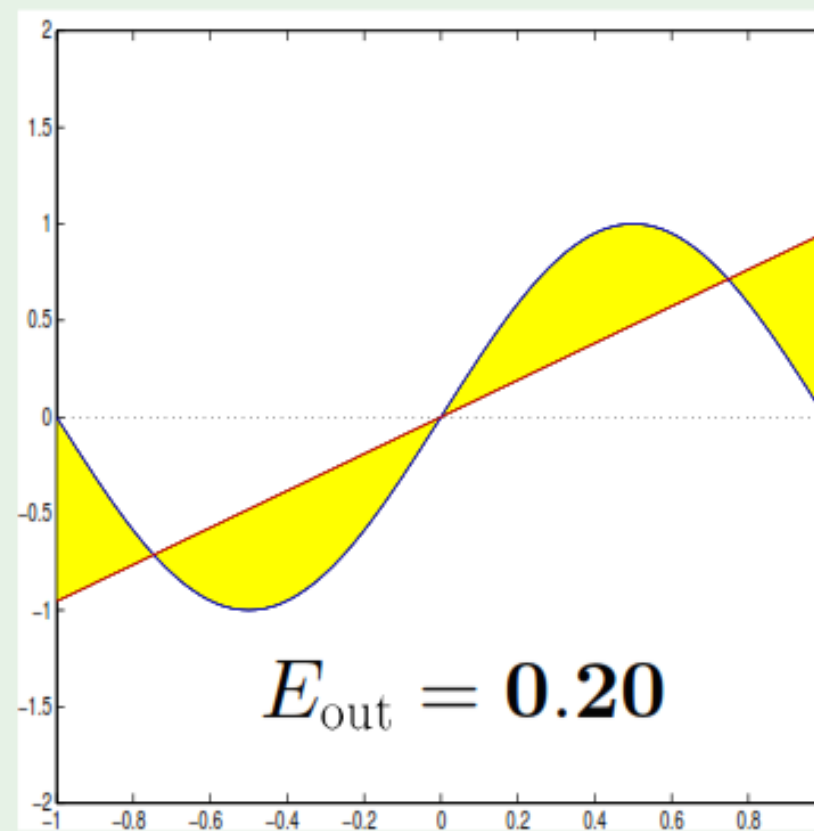


Approximation - \mathcal{H}_0 versus \mathcal{H}_1

\mathcal{H}_0

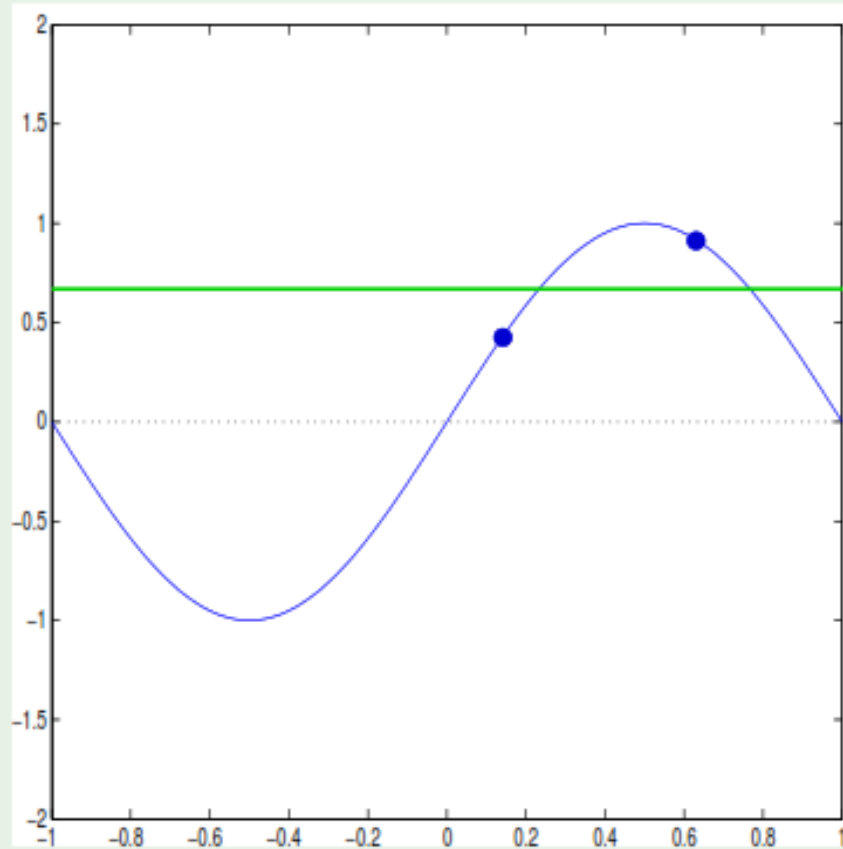


\mathcal{H}_1

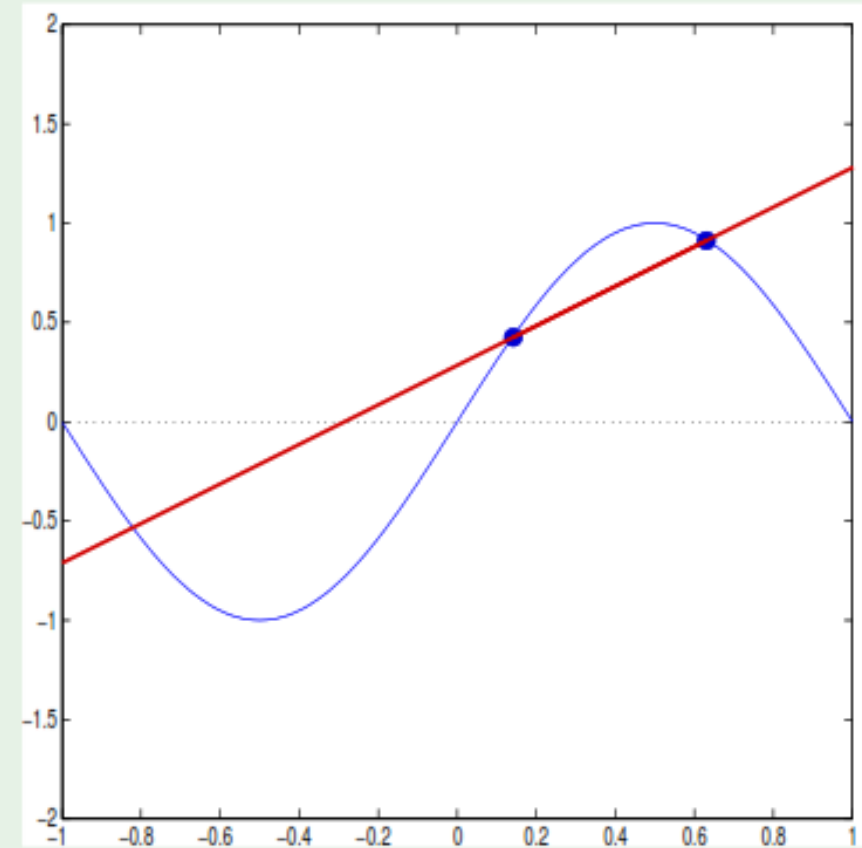


Learning - \mathcal{H}_0 versus \mathcal{H}_1

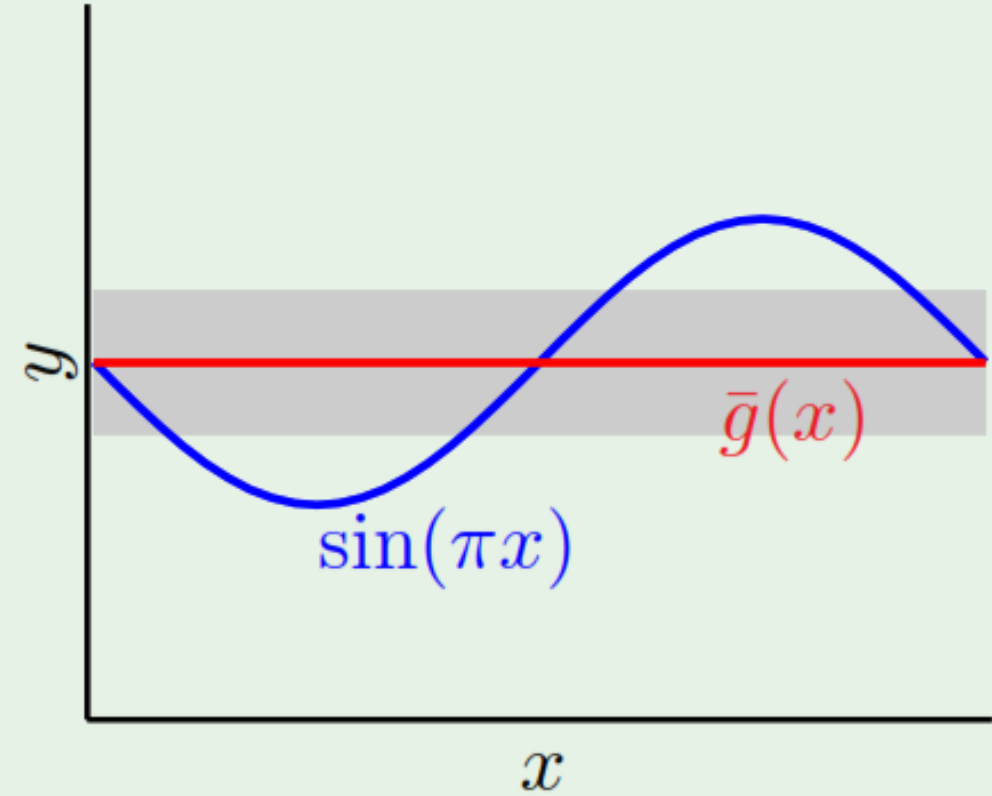
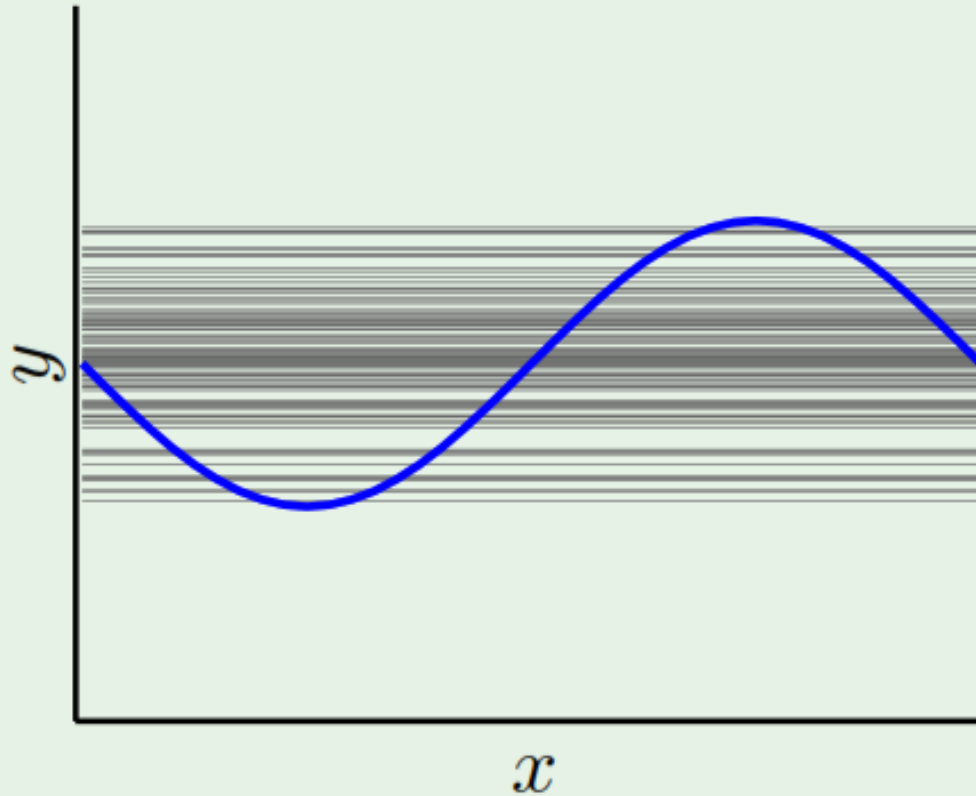
\mathcal{H}_0



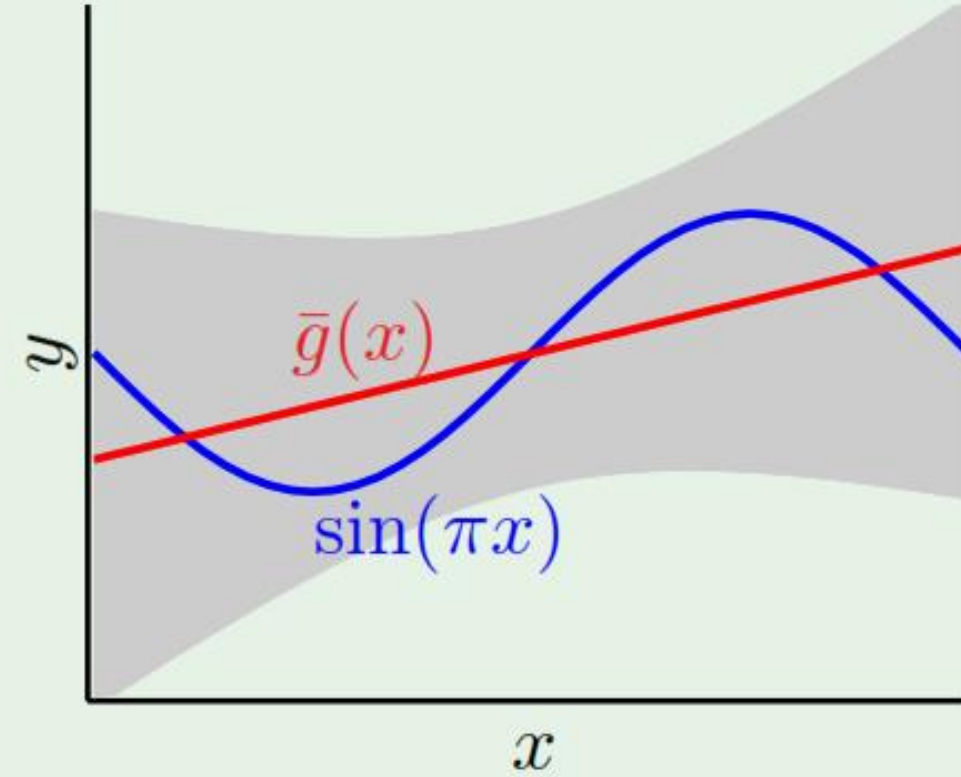
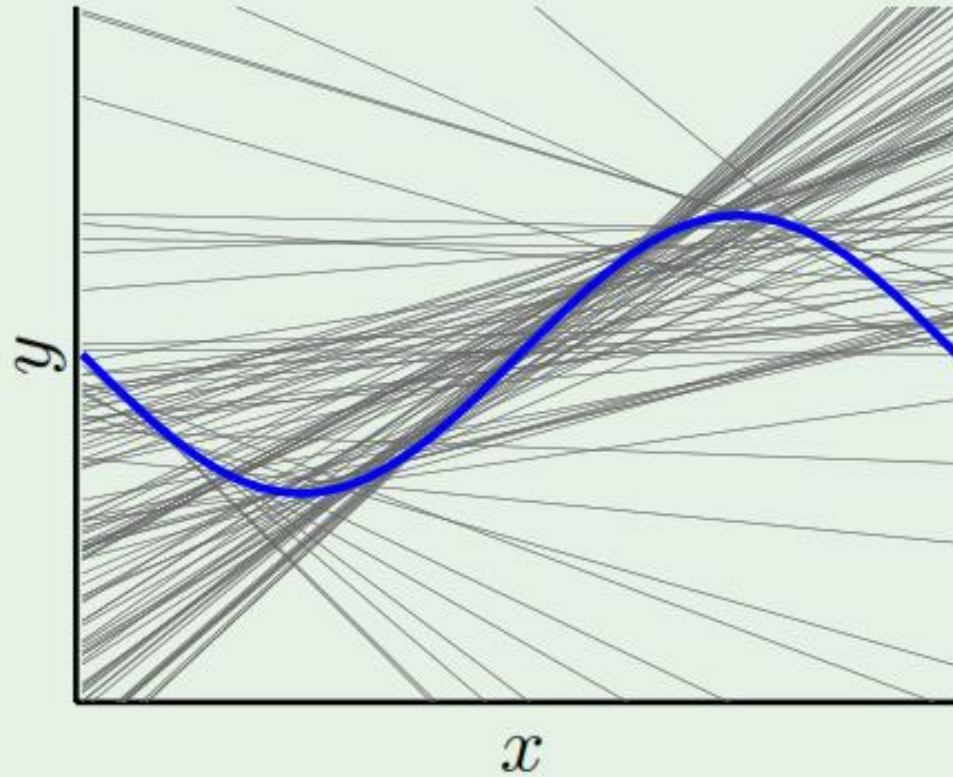
\mathcal{H}_1



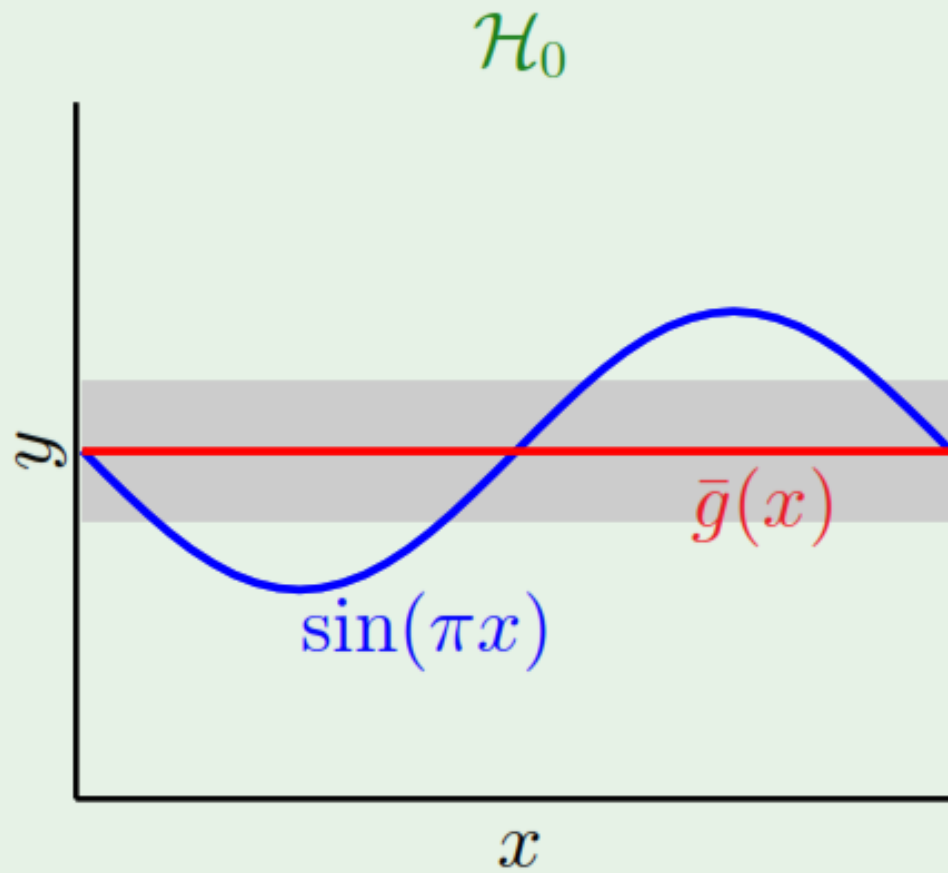
Bias and variance - \mathcal{H}_0



Bias and variance - \mathcal{H}_1

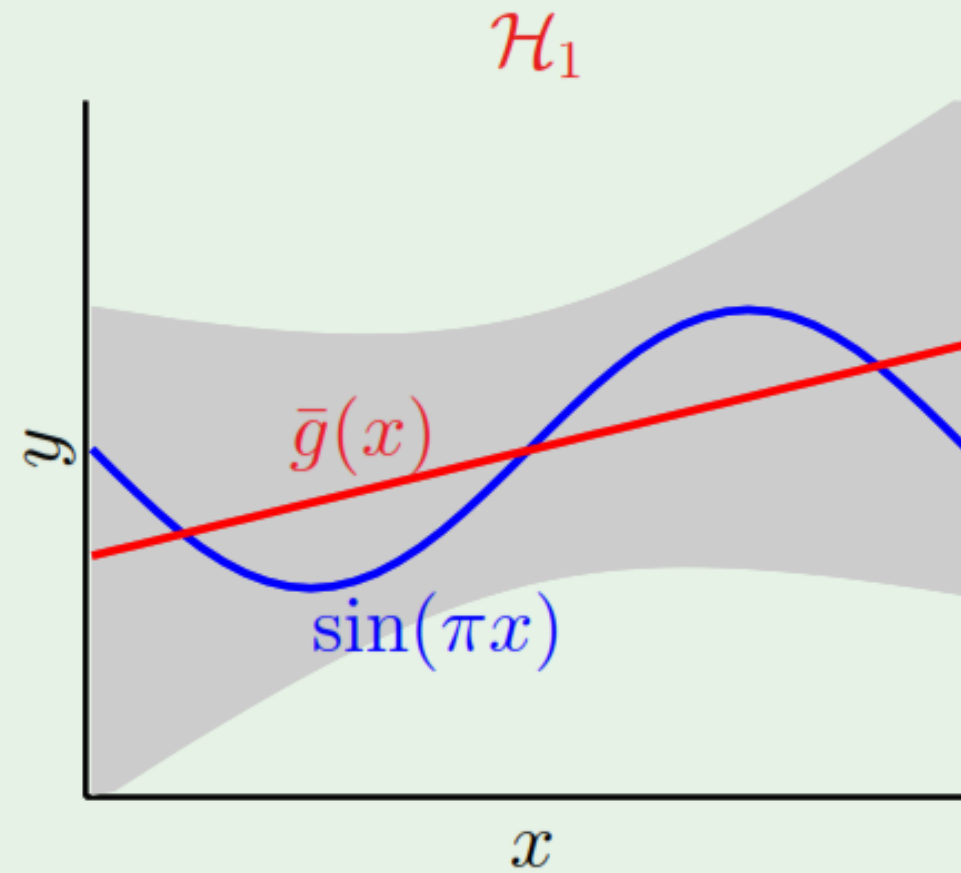


and the winner is ...



bias = **0.50**

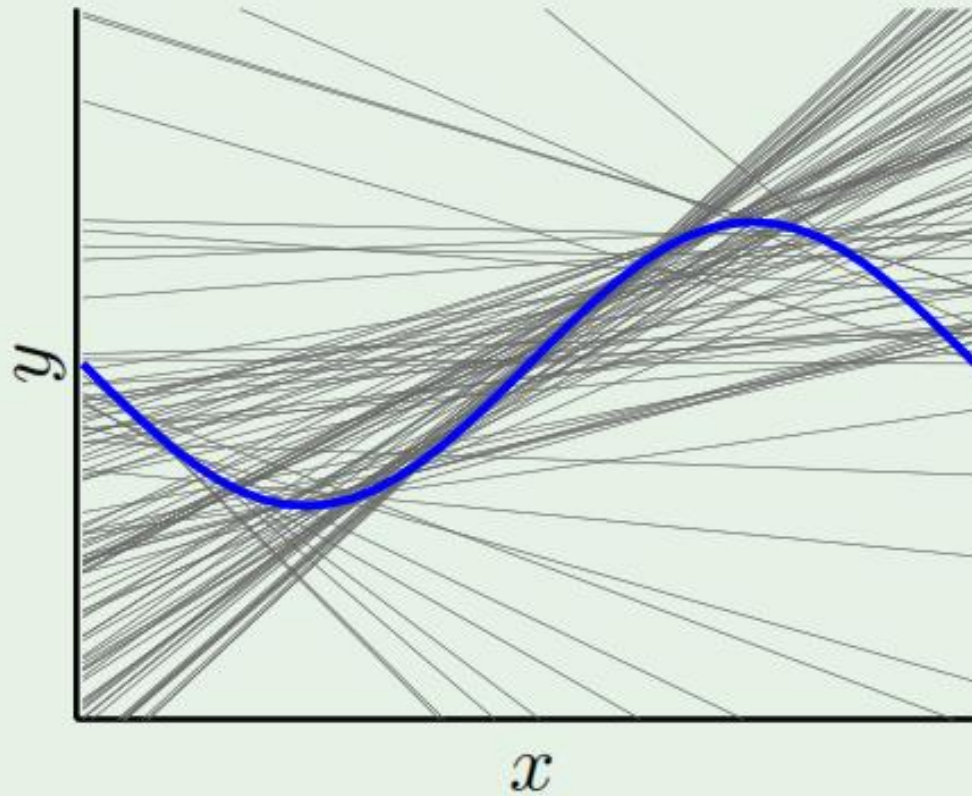
var = **0.25**



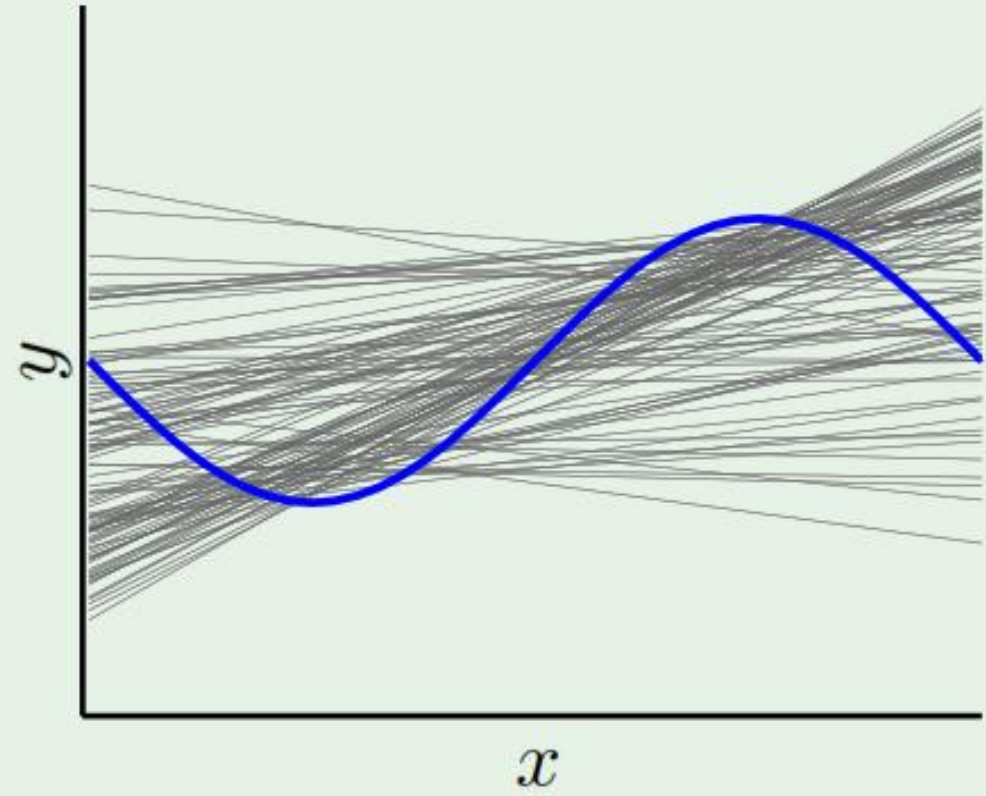
bias = **0.21**

var = **1.69**

A familiar example



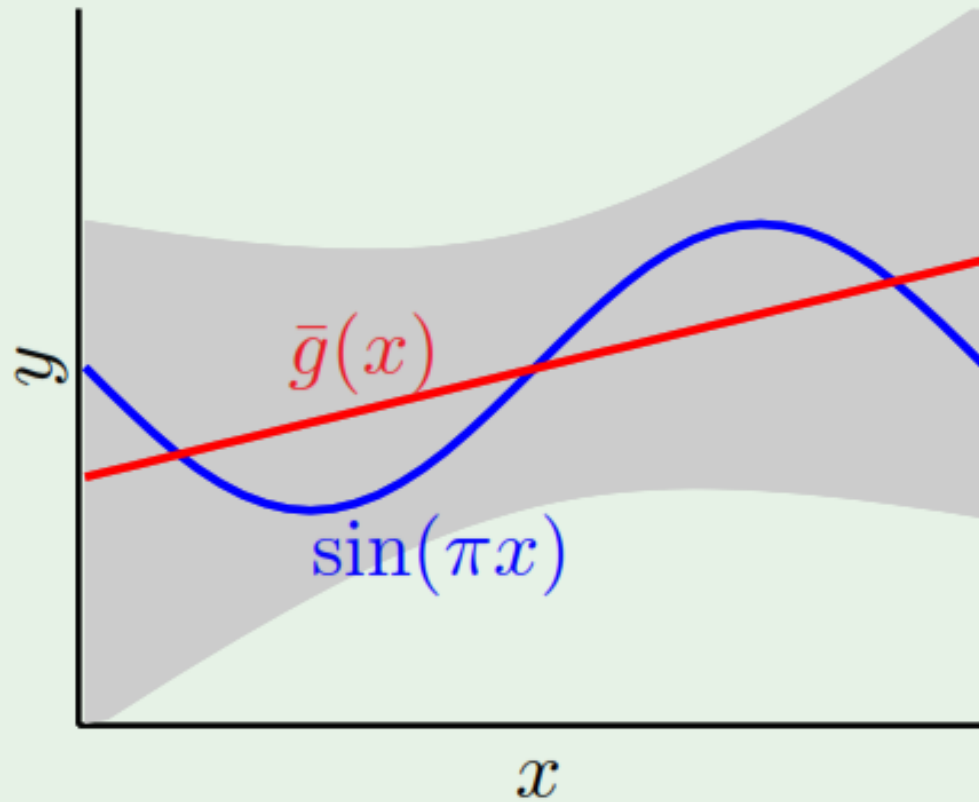
without regularization



with regularization

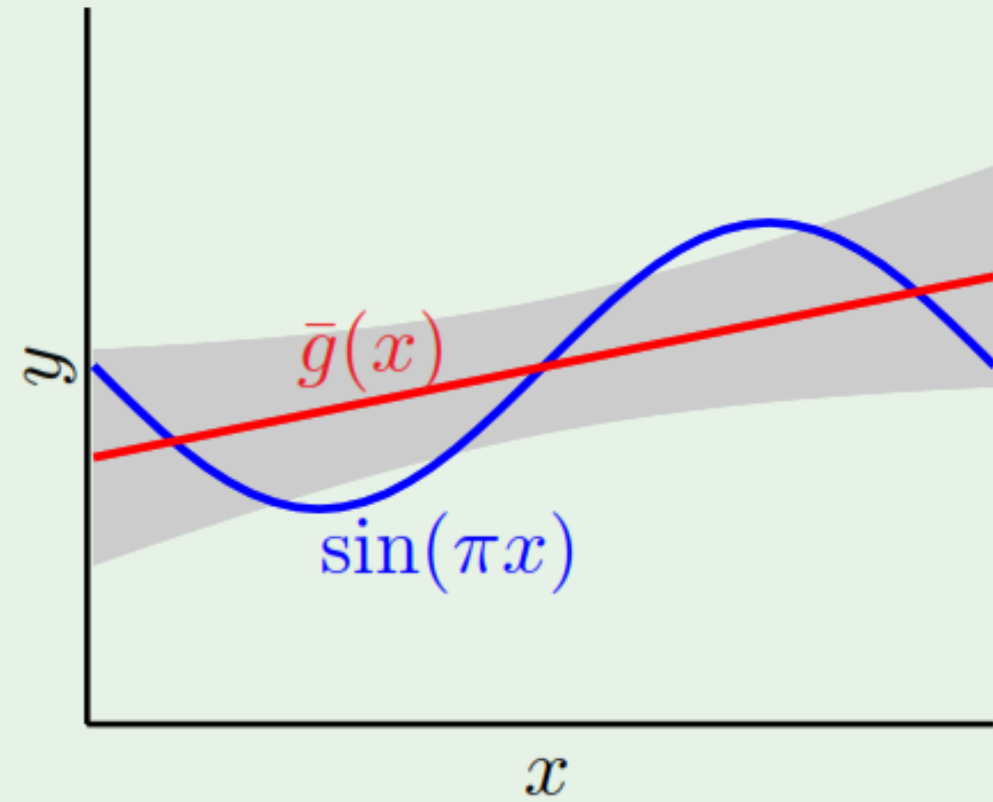
and the winner is ...

without regularization

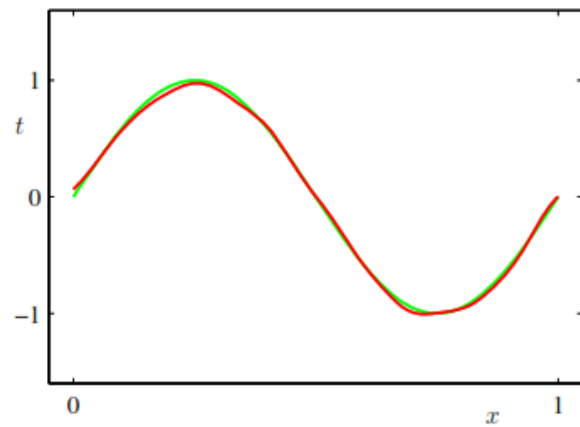
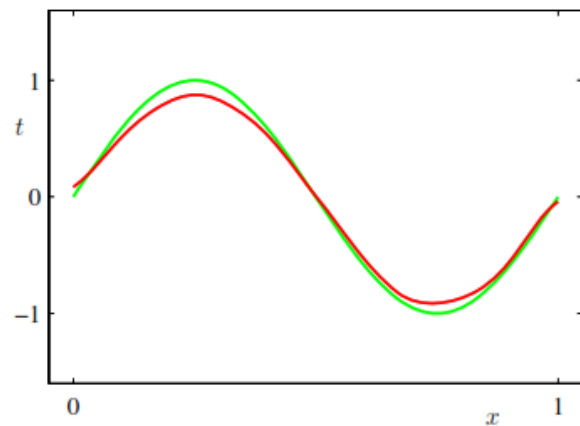
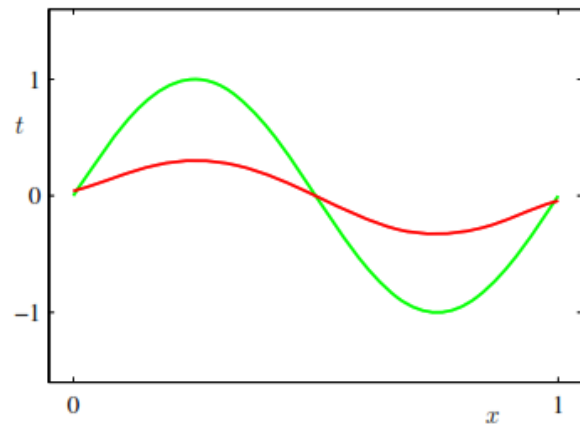


bias = **0.21** var = **1.69**

with regularization

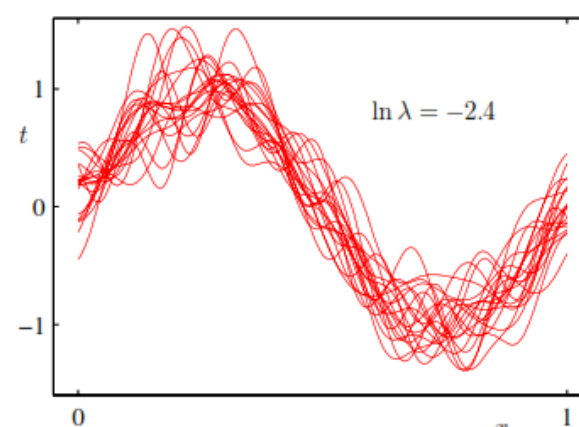
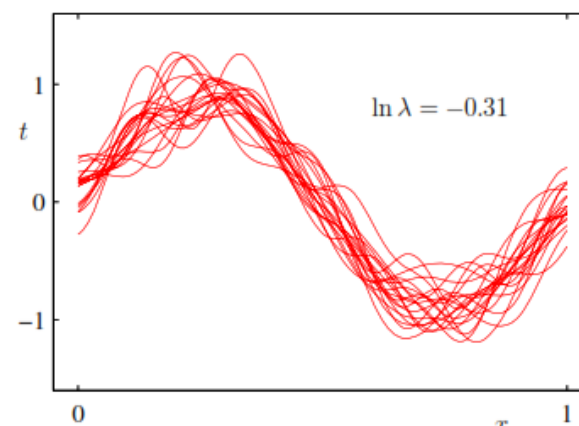
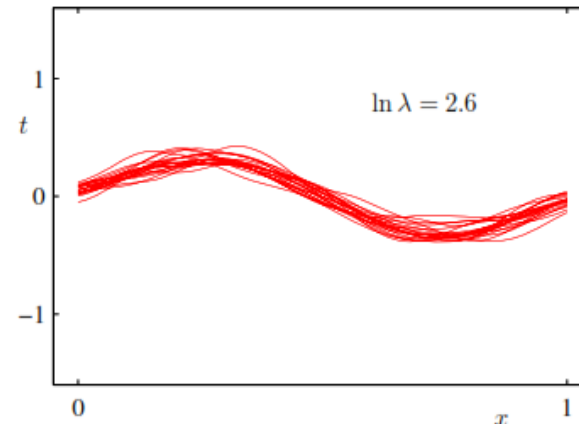


bias = **0.23** var = **0.33**



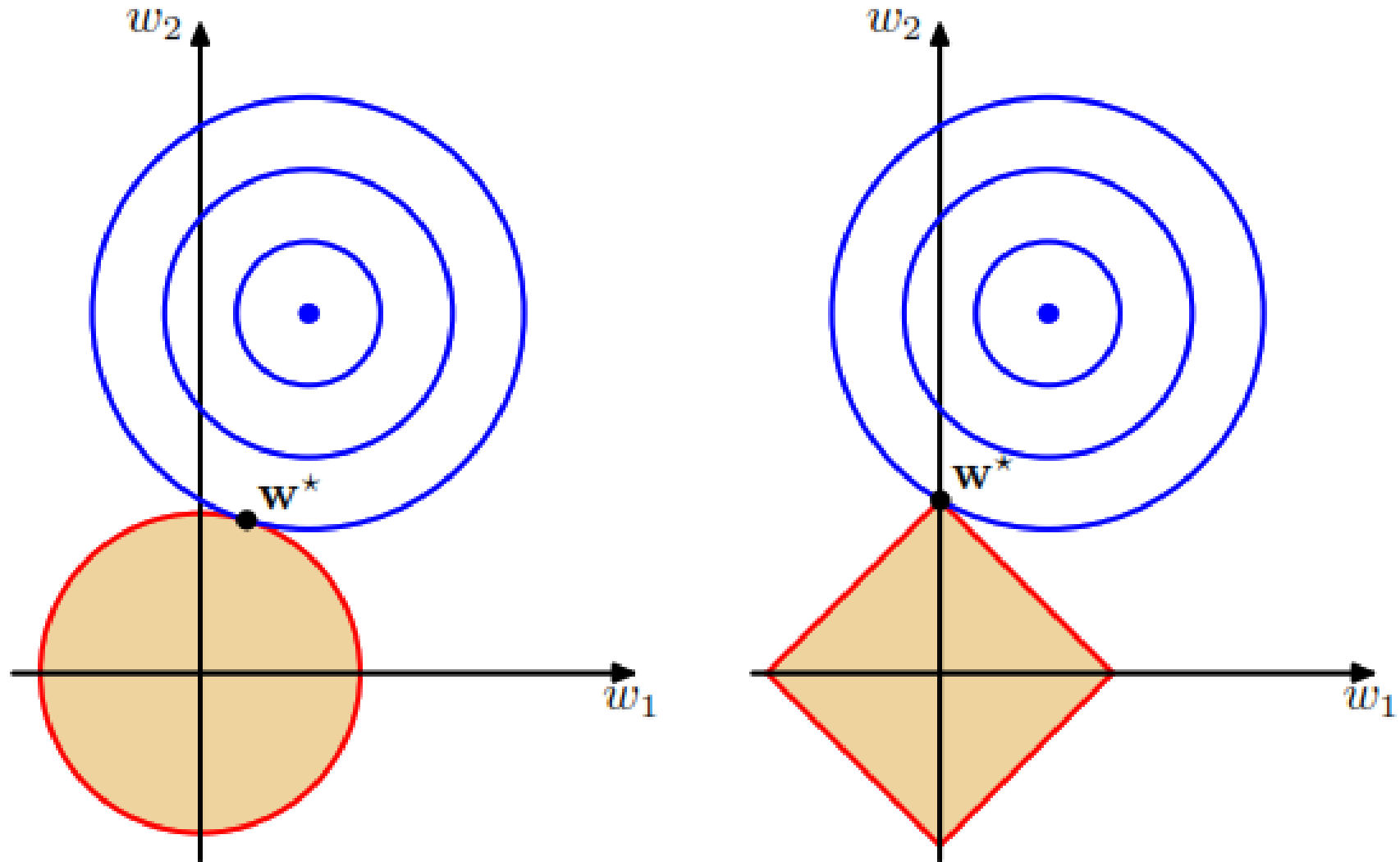
Average of 100 fits

100 datasets
 $n=100$
 $P=25$



Results of all fits

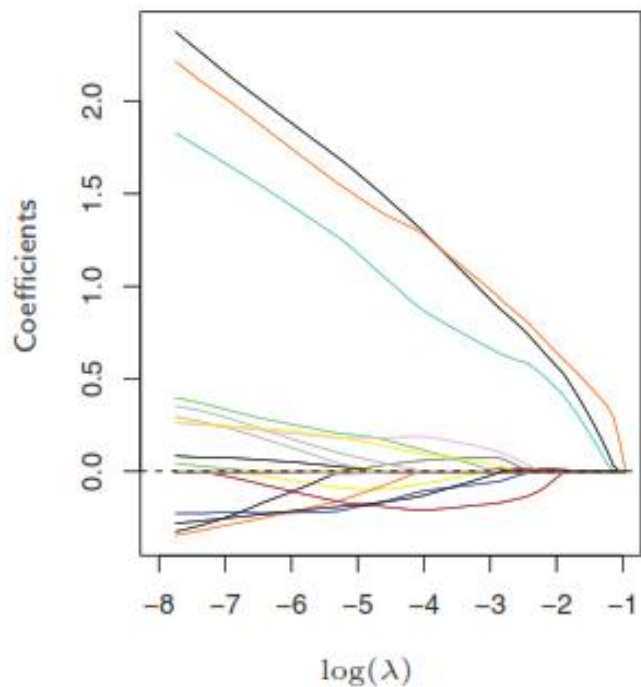
Contours of unregularized error function



Elastic net

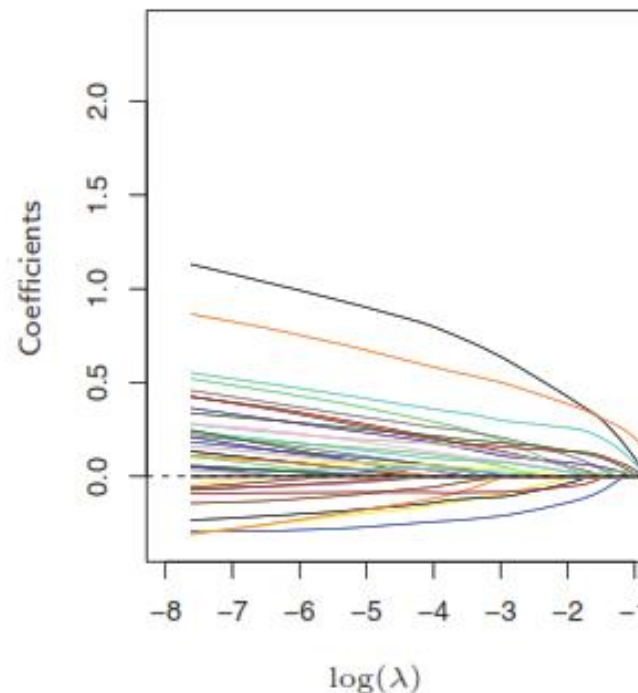
$$\hat{\beta} = \arg \min_{\beta} (Y - X\beta)^T (Y - X\beta) + \lambda(\alpha \|\beta\|_1 + (1 - \alpha) \|\beta\|_2^2)$$

Lasso



19 non-zero coefficients

Elastic Net



39 non-zero coefficients, but
with smaller magnitudes

References

- Pattern Recognition and Machine Learning by Christopher Bishop
- Learning from data by Abu-Mostafa, Y.S., Magdon-Ismael, M. and Lin, H.T.
 - Slides 7-14 are from the lectures 8 and 12 of *Learning from data* course at Caltech
 - <https://work.caltech.edu/lectures.html>