

Lecture 5

Machine Learning

Last Times:

- Expectations, sample average
- The Law of large numbers and Monte Carlo
- Sampling Methods

Law of Large numbers (LLN)

- Expectations become sample averages. Convergence for large N.

$$\begin{aligned} E_f[g] &= \int g(x) dF = \int g(x) f(x) dx \\ &= \lim_{n \rightarrow \infty} \frac{1}{N} \sum_{x_i \sim f} g(x_i) \end{aligned}$$

- for finite N a sample average
- thus expectations in the replication "dimension" come into play
- mean of sample means and standard error
- this is the sampling distribution
- CLT and all that jazz

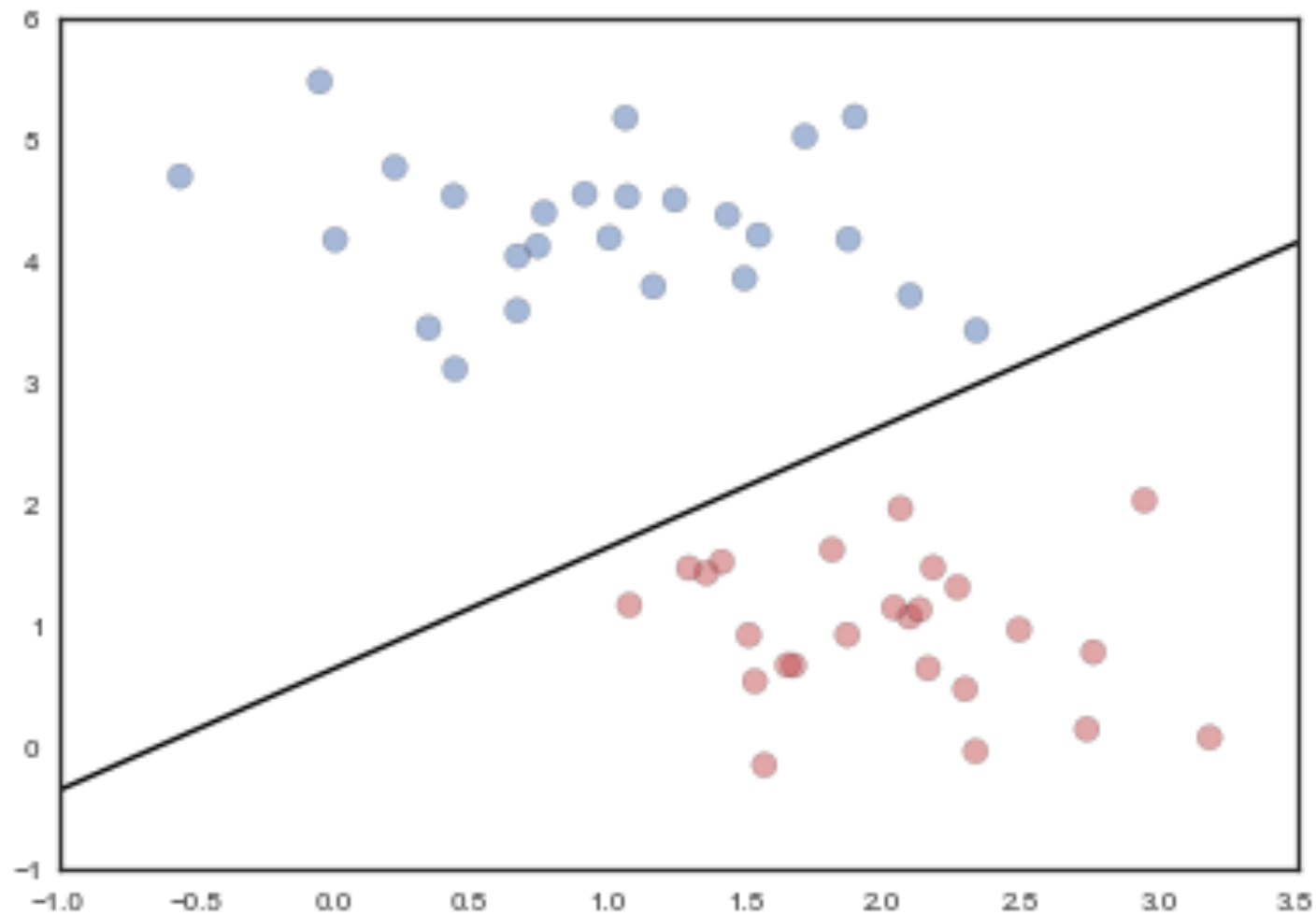
Today: machine Learning

- noiseless models, the approximation problems
- models with noise
- test sets and learning theory
- validation and cross-validation
- regularization

Why study this?

- isn't this a course in Stoch Opt and Bayes?
- application of law of large numbers
- establishes ideas of supervised learning
- learn validation for model selection
- bayes critical to understand machine learning

CLASSIFICATION

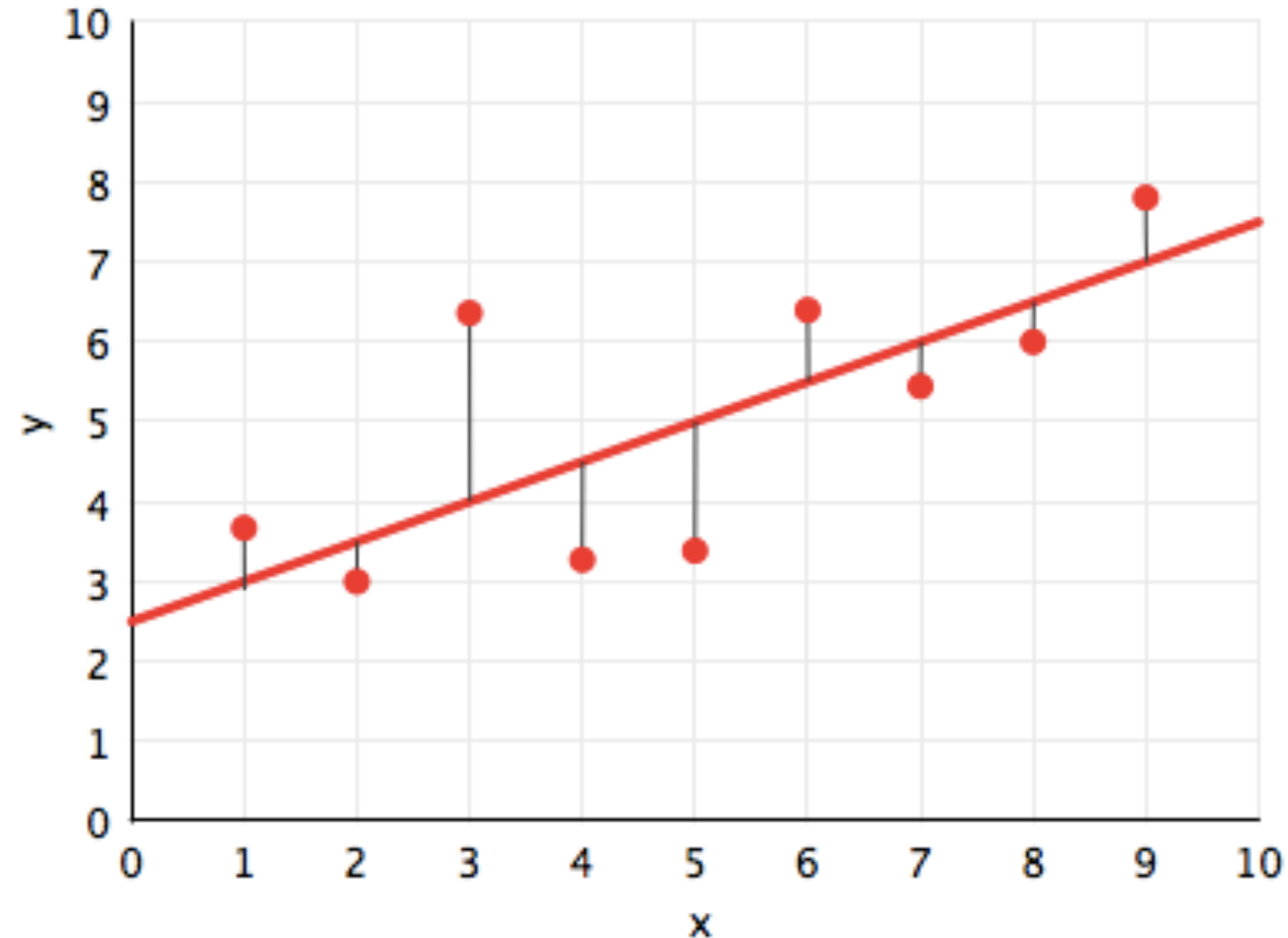


- will a customer churn?
- is this a check? For how much?
- a man or a woman?
- will this customer buy?
- do you have cancer?
- is this spam?
- whose picture is this?

^j image from code in <http://bit.ly/1Azg29G>

REGRESSION

- how many dollars will you spend?
- what is your creditworthiness
- how many people will vote for Bernie t days before election
- use to predict probabilities for classification
- causal modeling in econometrics



From Bayesian Reasoning and Machine Learning, David Barber:

"A father decides to teach his young son what a sports car is. Finding it difficult to explain in words, he decides to give some examples. They stand on a motorway bridge and ... the father cries out 'that's a sports car!' when a sports car passes by. After ten minutes, the father asks his son if he's understood what a sports car is. The son says, 'sure, it's easy'. An old red VW Beetle passes by, and the son shouts – 'that's a sports car!'. Dejected, the father asks – 'why do you say that?'. 'Because all sports cars are red!', replies the son."

HYPOTHESIS SPACES

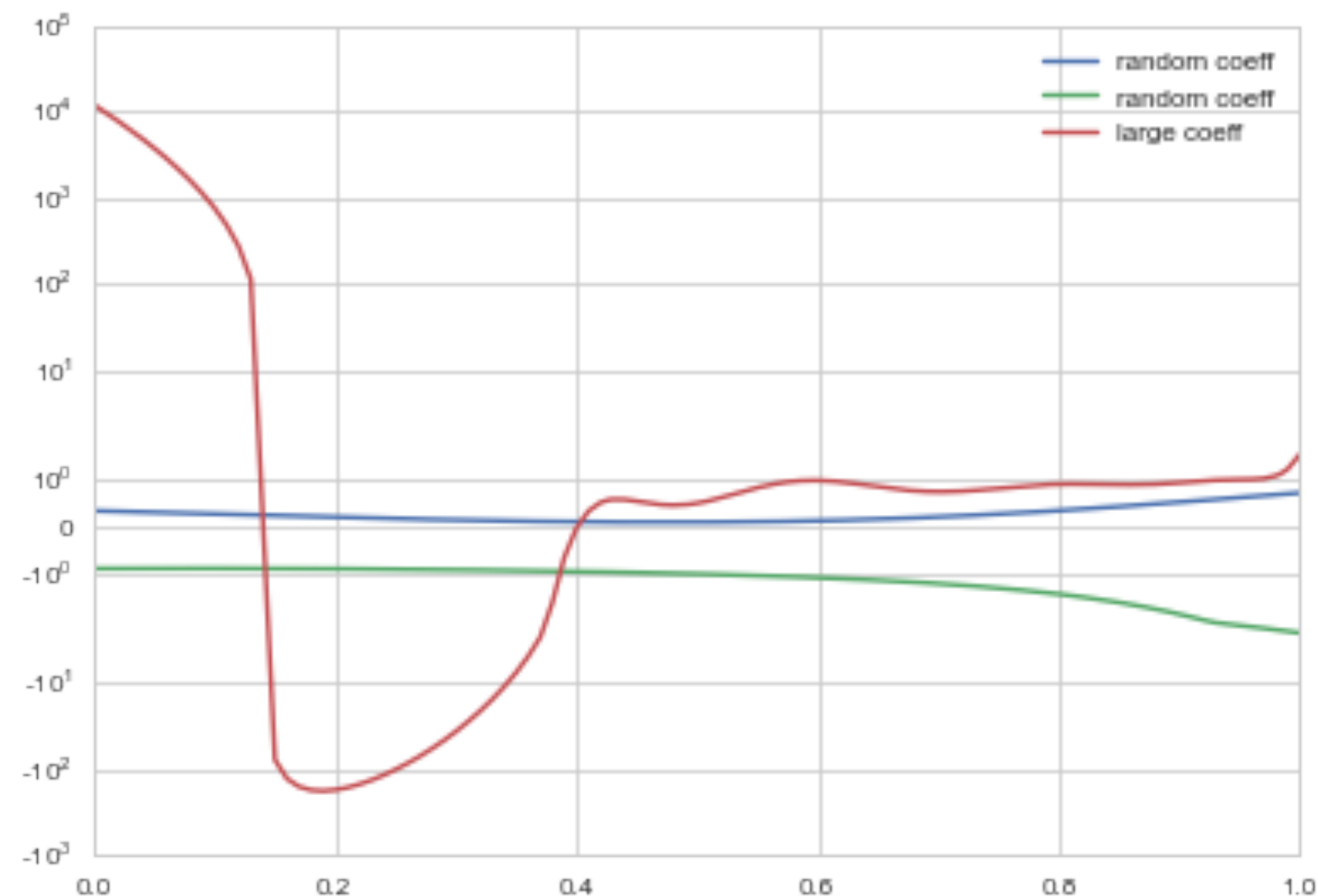
A polynomial looks so:

$$h(x) = \theta_0 + \theta_1 x^1 + \theta_2 x^2 + \dots + \theta_n x^n = \sum_{i=0}^n \theta_i x^i$$

All polynomials of a degree or complexity d constitute a hypothesis space.

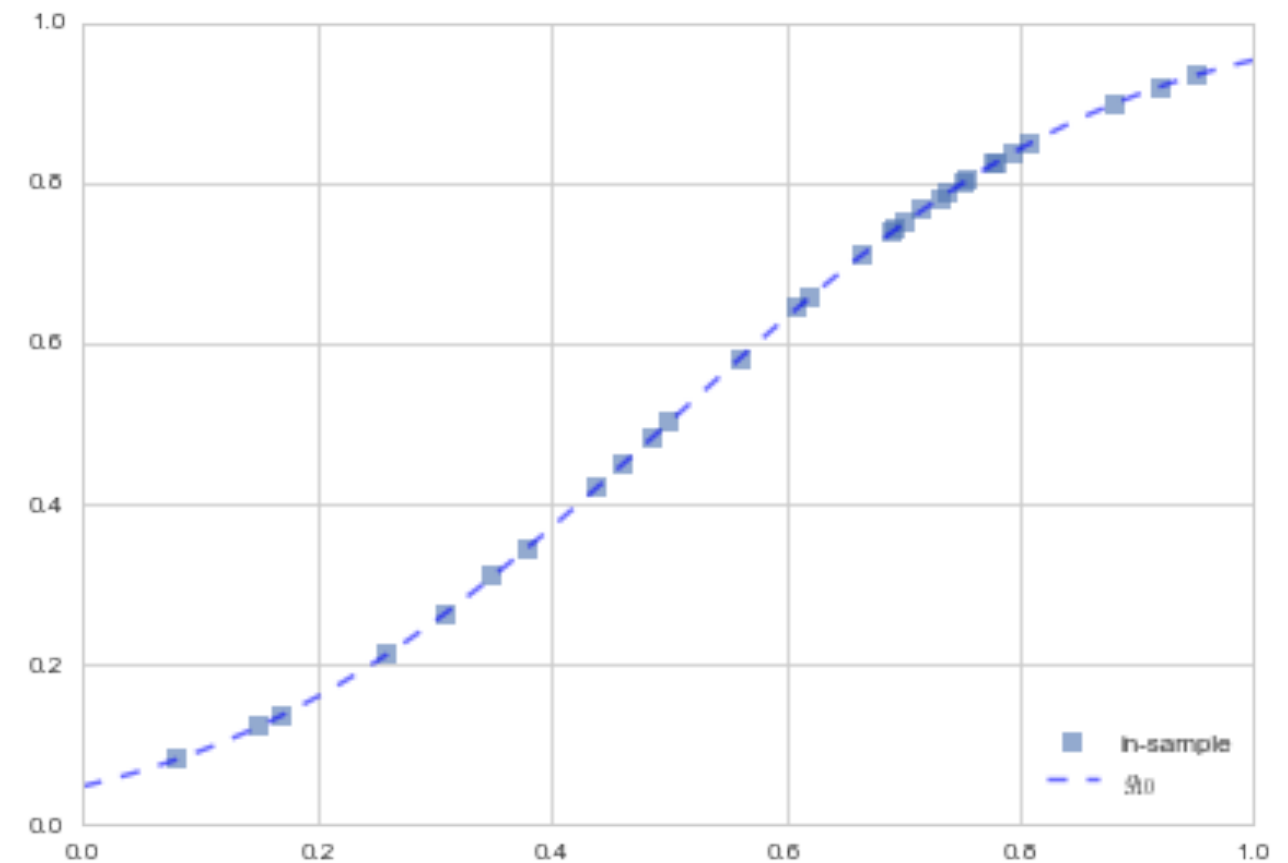
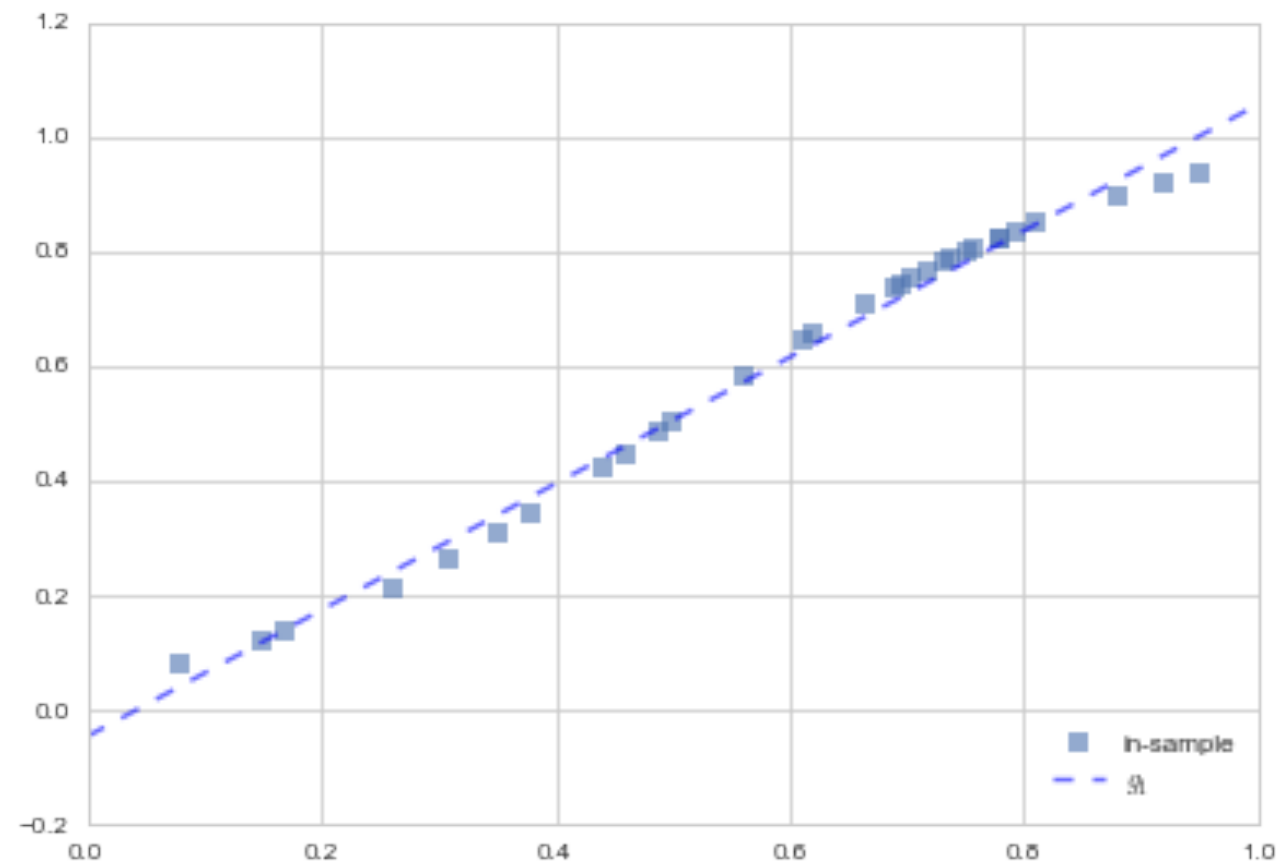
$$\mathcal{H}_1 : h_1(x) = \theta_0 + \theta_1 x$$

$$\mathcal{H}_{20} : h_{20}(x) = \sum_{i=0}^{20} \theta_i x^i$$

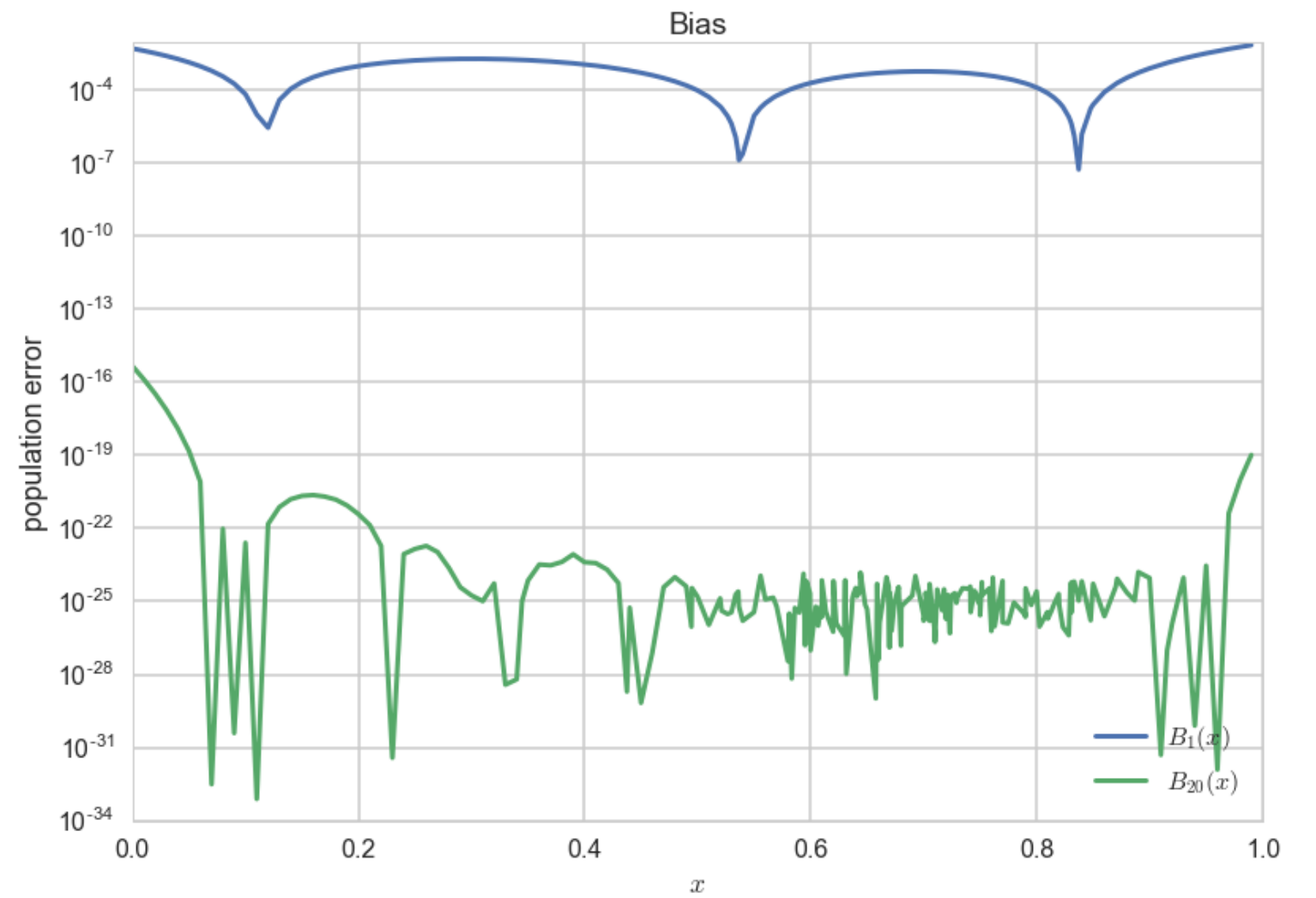
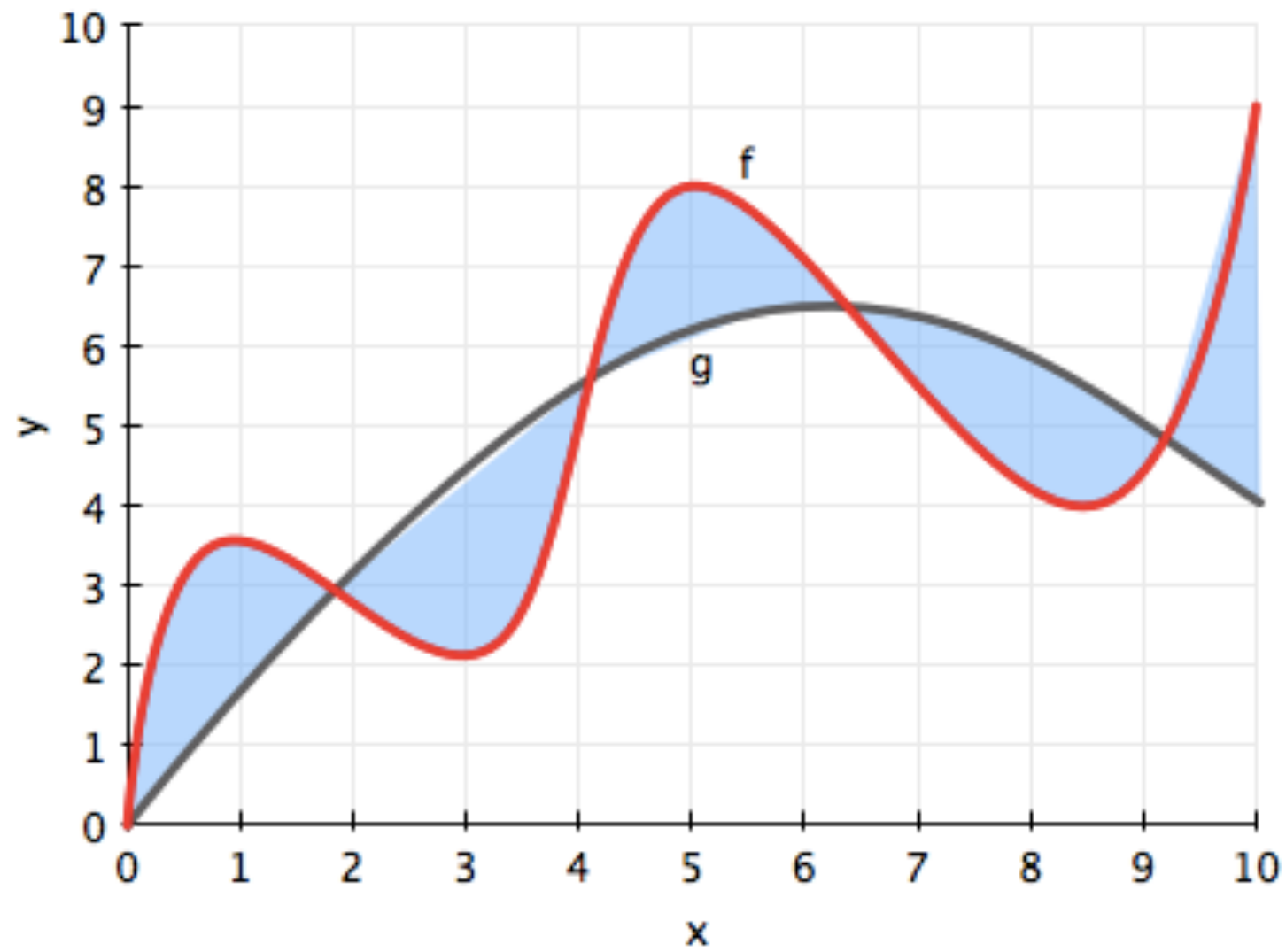


Approximation: Learning without noise

30 points of data. Which fit is better? Line in \mathcal{H}_1 or curve in \mathcal{H}_{20} ?



Bias



RISK: What does it mean to FIT?

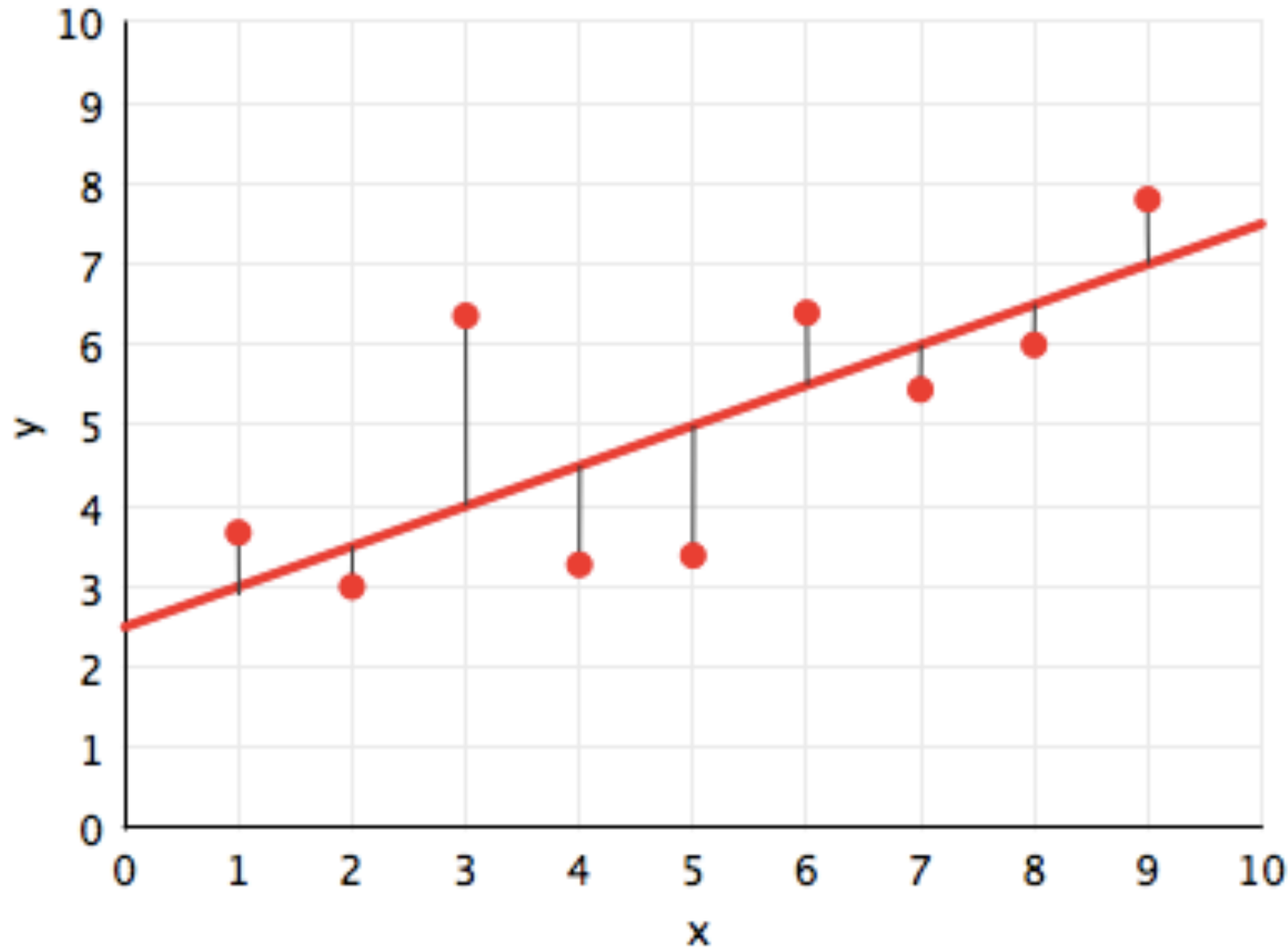
Minimize distance from the line?

$$R_{\mathcal{D}}(h_1(x)) = \frac{1}{N} \sum_{y_i \in \mathcal{D}} (y_i - h_1(x_i))^2$$

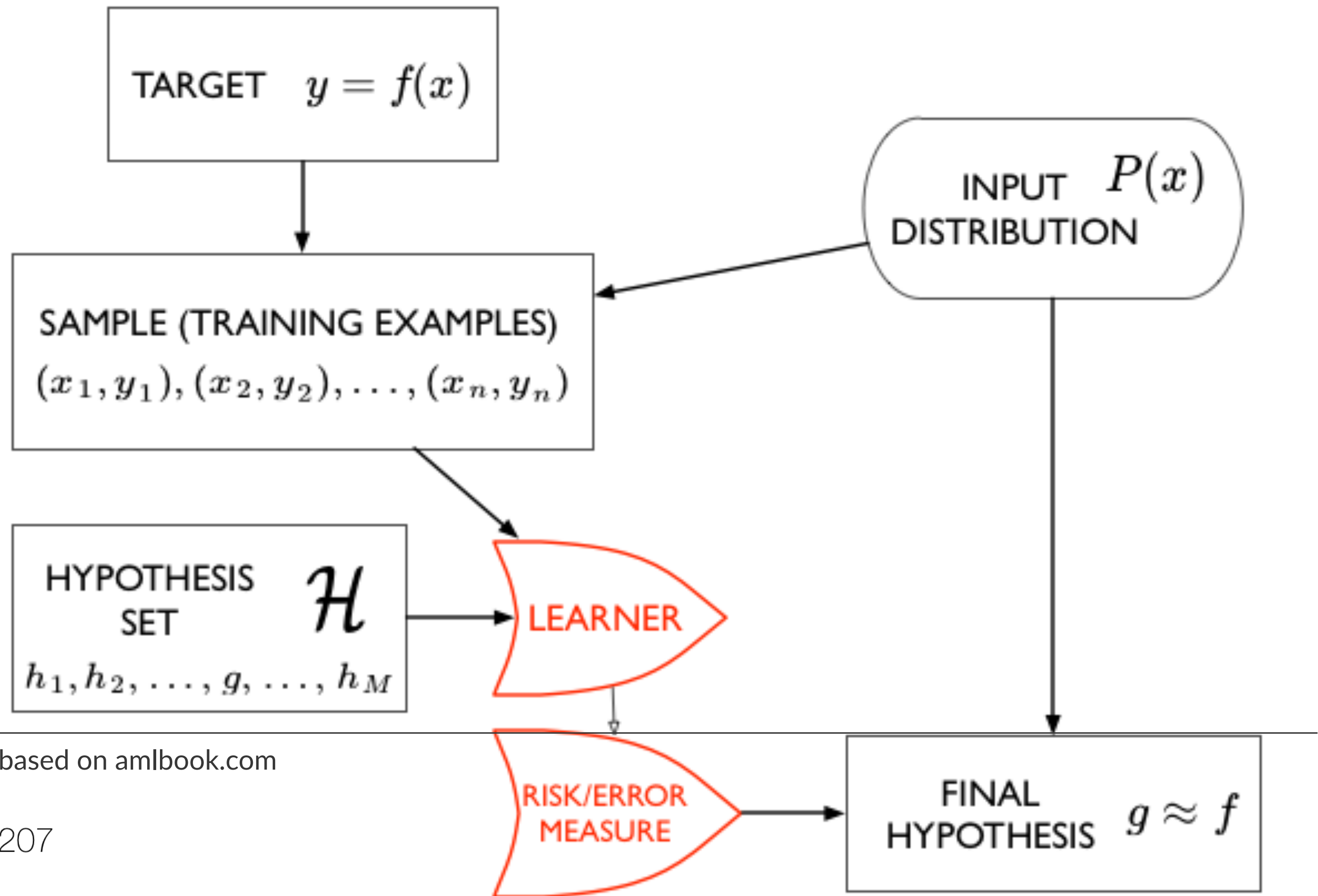
Minimize squared distance from the line.
Empirical Risk Minimization.

$$g_1(x) = \arg \min_{h_1(x) \in \mathcal{H}} R_{\mathcal{D}}(h_1(x)).$$

Get intercept w_0 and slope w_1 .



*



* image based on amlbook.com

SAMPLE vs POPULATION

$$\text{Want: } R_{out}(h) = E_{p(x)} [(h(x) - f(x))^2] = \int dx p(x) (h(x) - f(x))^2$$

LLN:

$$R_{out}(h) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{x_i \sim p(x)} (h(x_i) - f(x_i))^2 = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{x_i \sim p(x)} (h(x_i) - y_i)^2$$

$$\mathcal{D} \text{ representative } (\mathcal{D} \sim p(x)) \implies \mathcal{R}_{\mathcal{D}}(h) = \sum_{x_i \in \mathcal{D}} (h(x_i) - y_i)^2$$

Statement of the Learning Problem

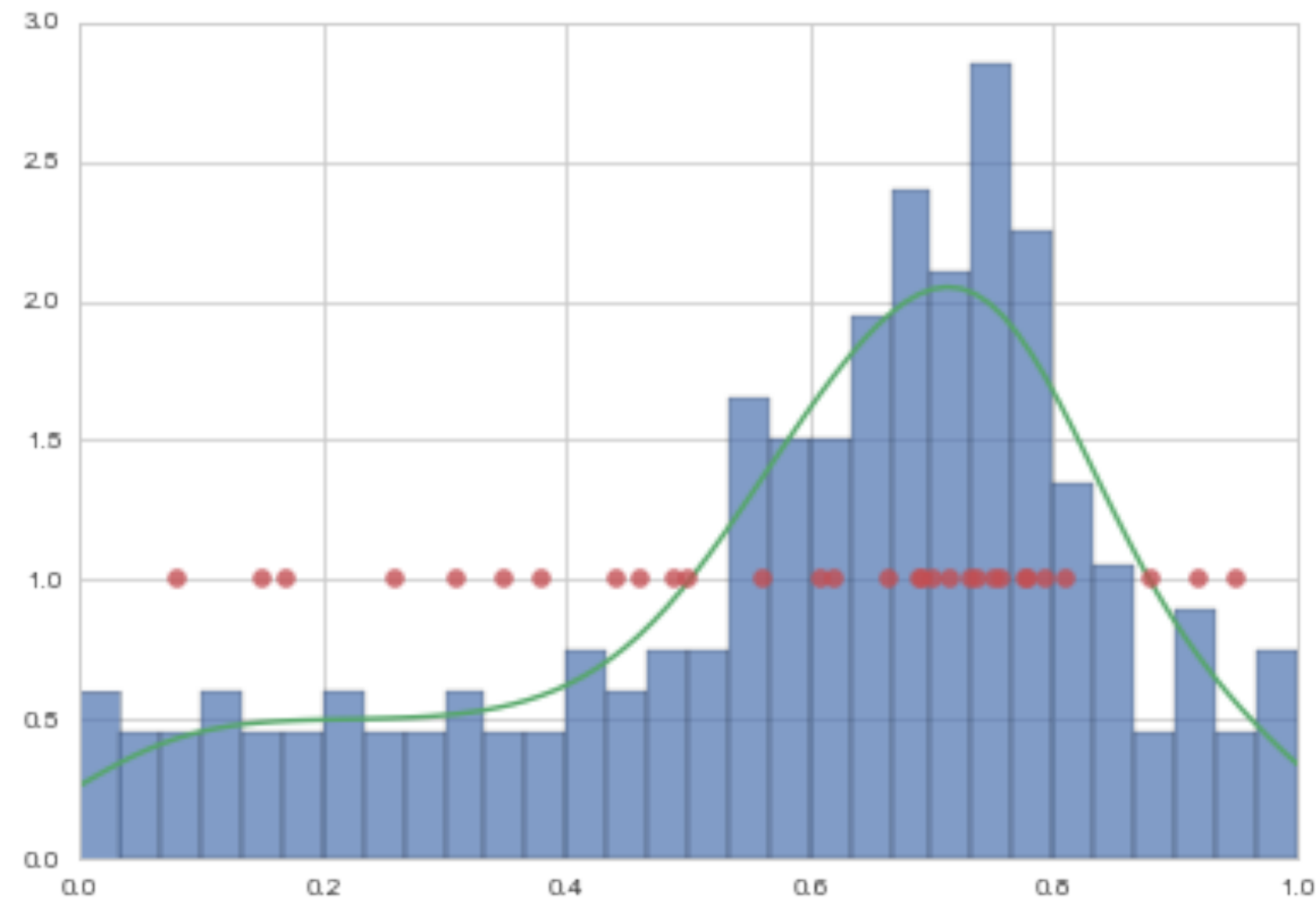
The sample must be representative of the population!

$A : R_{\mathcal{D}}(g)$ *smallest on \mathcal{H}*

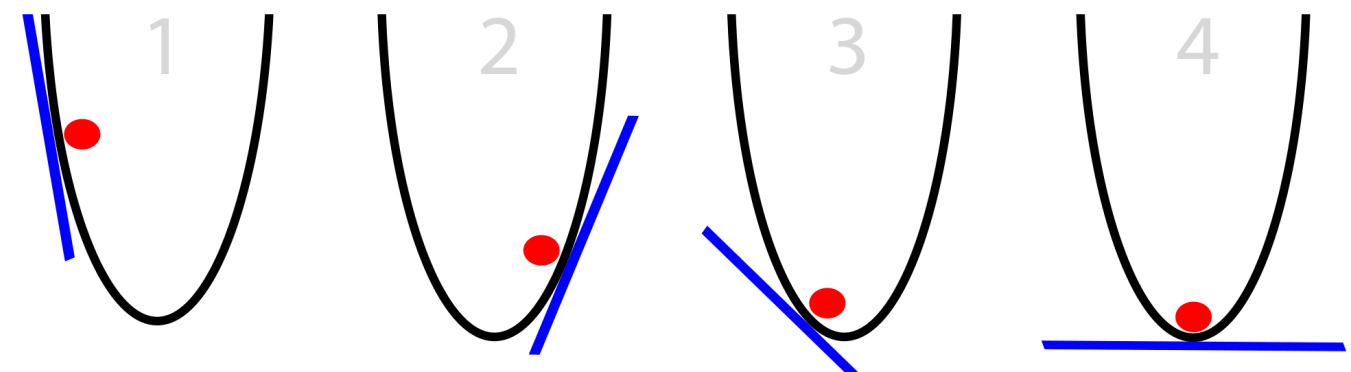
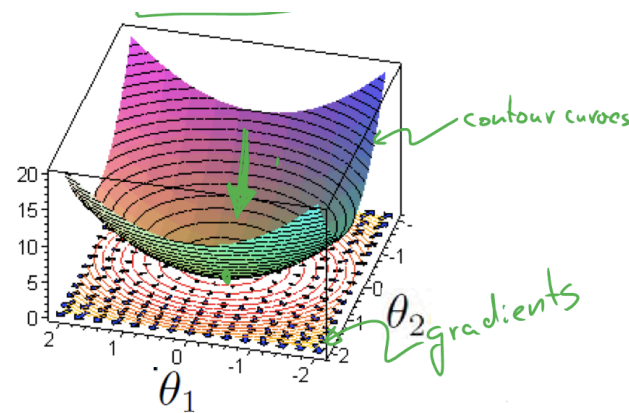
$B : R_{out}(g) \approx R_{\mathcal{D}}(g)$

A: Empirical risk estimates in-sample risk.

B: Thus the out of sample risk is also small.



CONVEX MINIMIZATION

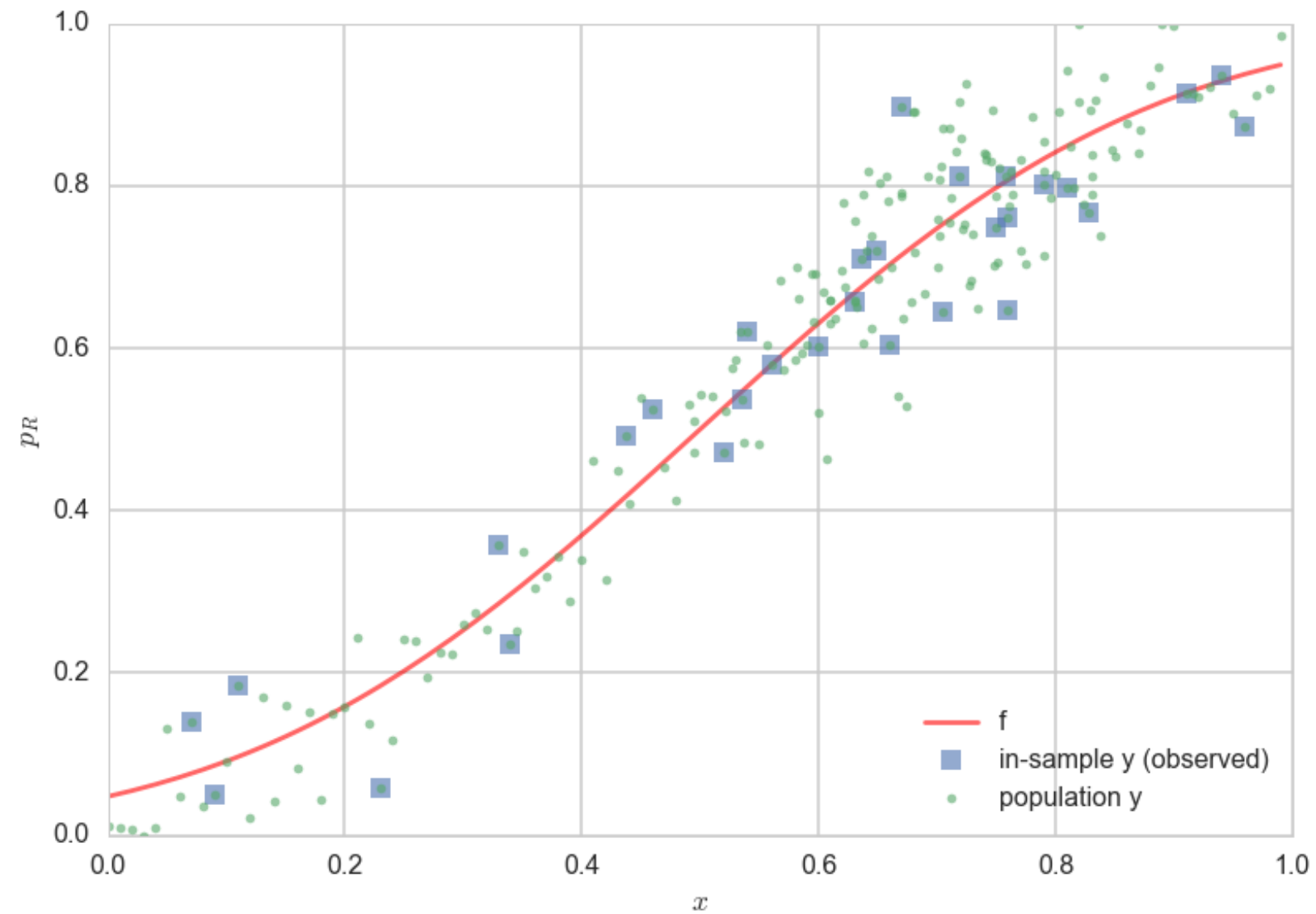


In general one can use gradient descent .

For linear-regression, one can however just do this using matrix algebra.

Image From Nando-deFreitas Deep Learning Course 2015

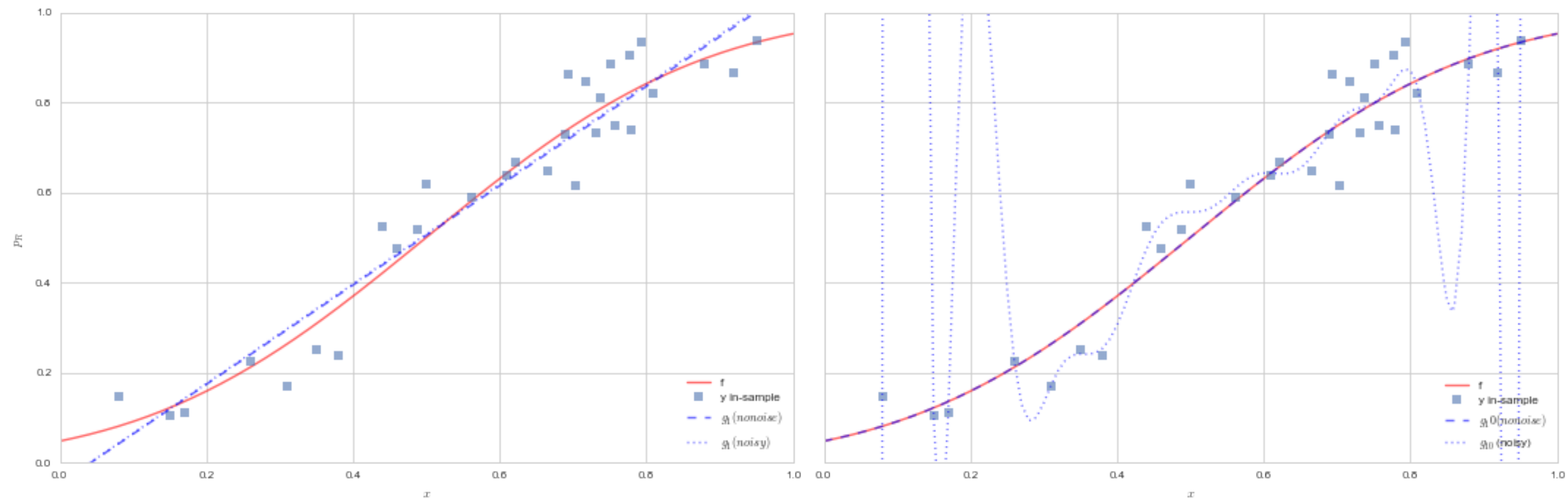
THE REAL WORLD HAS NOISE



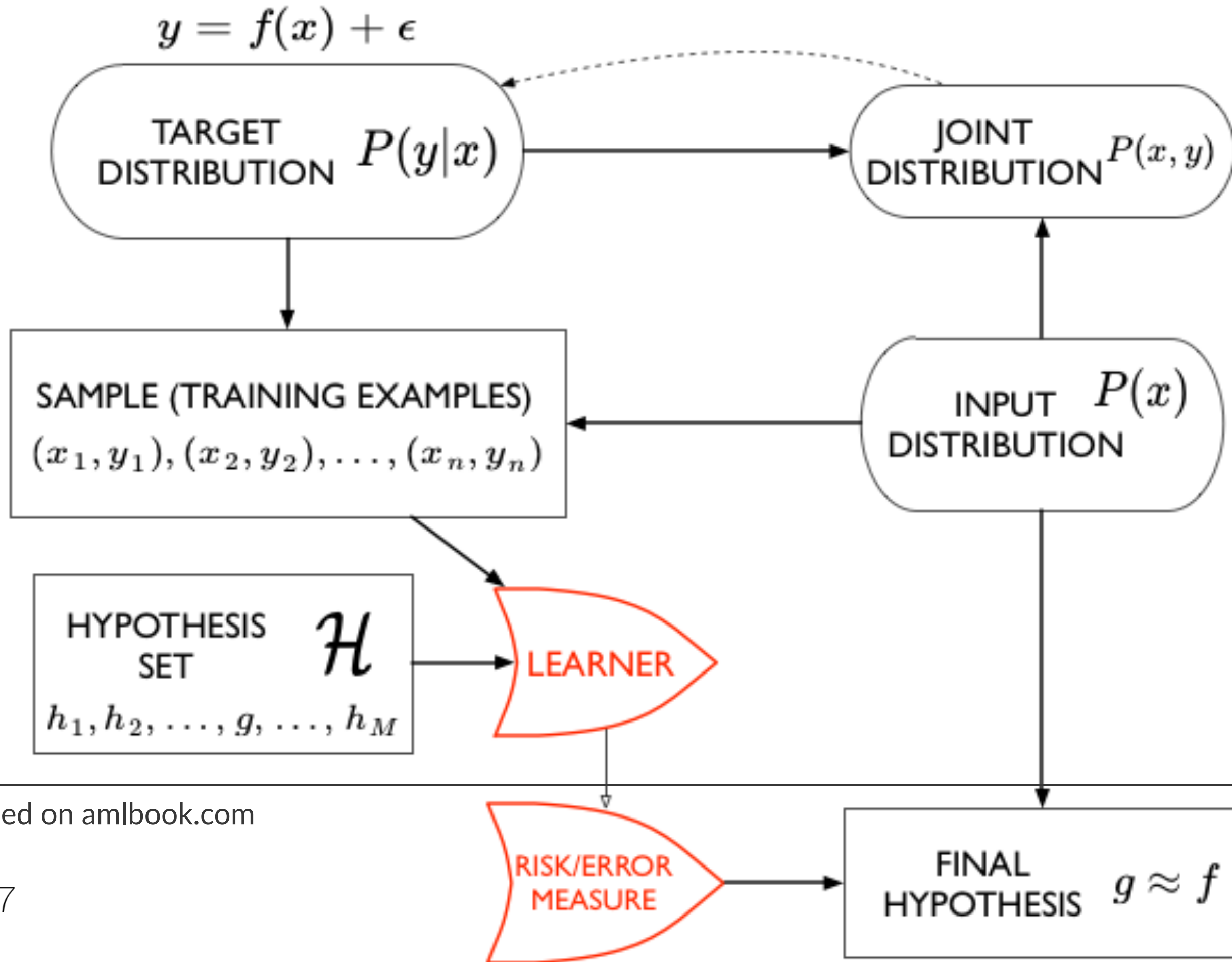
THE REAL WORLD HAS NOISE

Which fit is better now?

The line or the curve?

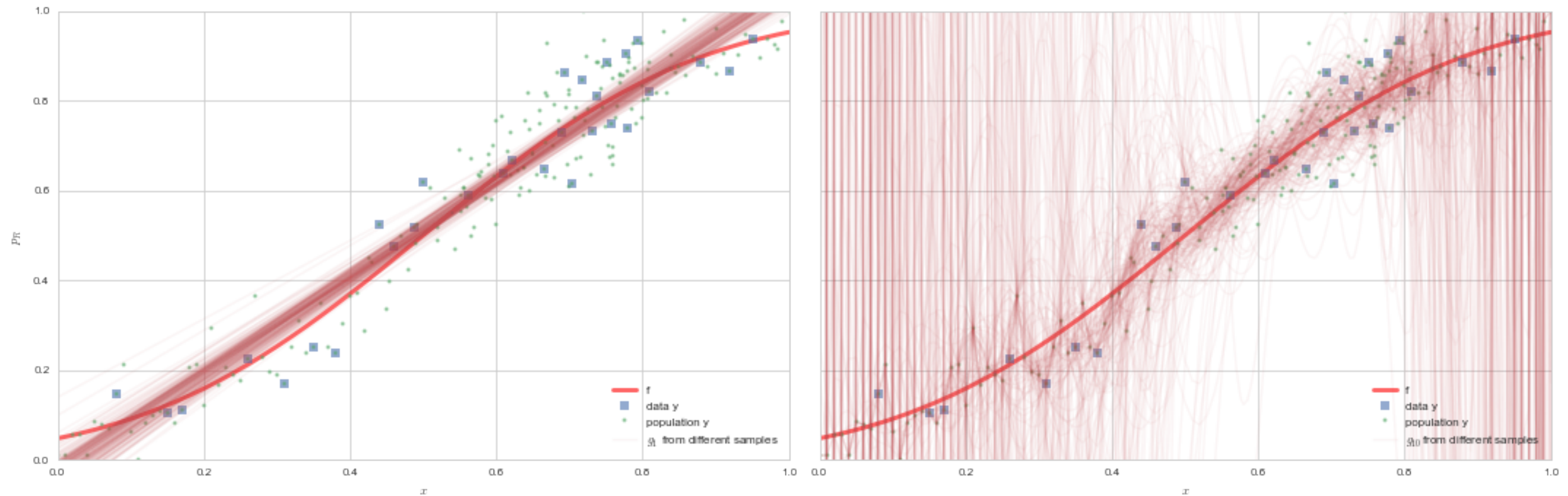


*



* image based on amlbook.com

UNDERFITTING (Bias) vs OVERFITTING (Variance)



Every model has Bias and Variance

$$R_{out}(h) = E_{p(x)}[(h(x) - y)^2] = \int dx p(x)(h(x) - f(x) - \epsilon)^2.$$

Fit hypothesis $h = g_{\mathcal{D}}$, where \mathcal{D} is our training sample.

Define:

$$\langle R \rangle = E_{\mathcal{D}}[R_{out}(g_{\mathcal{D}})] = E_{\mathcal{D}} E_{p(x)}[(g_{\mathcal{D}}(x) - f(x) - \epsilon)^2]$$

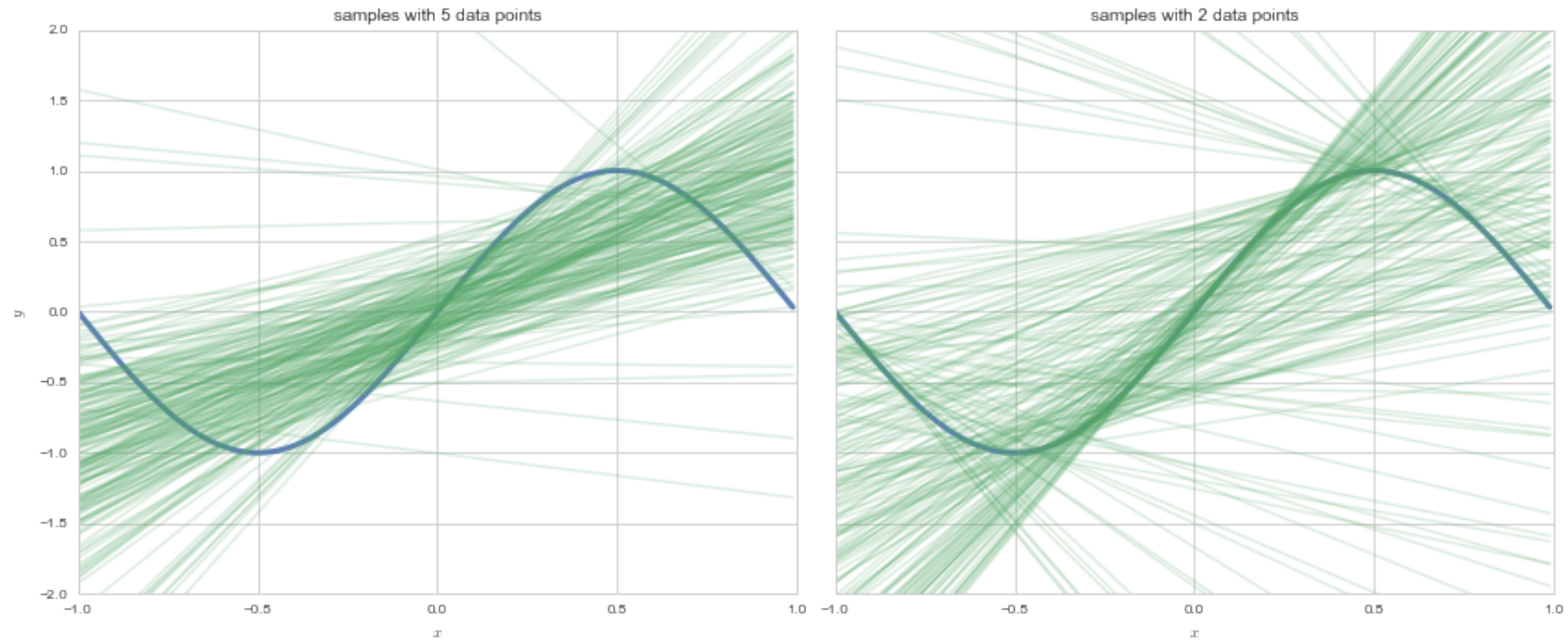
$$\bar{g} = E_{\mathcal{D}}[g_{\mathcal{D}}] = (1/M) \sum_{\mathcal{D}} g_{\mathcal{D}}$$

Then,

$$\langle R \rangle = E_{p(x)} [E_{\mathcal{D}} [(g_{\mathcal{D}} - \bar{g})^2]] + E_{p(x)} [(f - \bar{g})^2] + \sigma^2$$

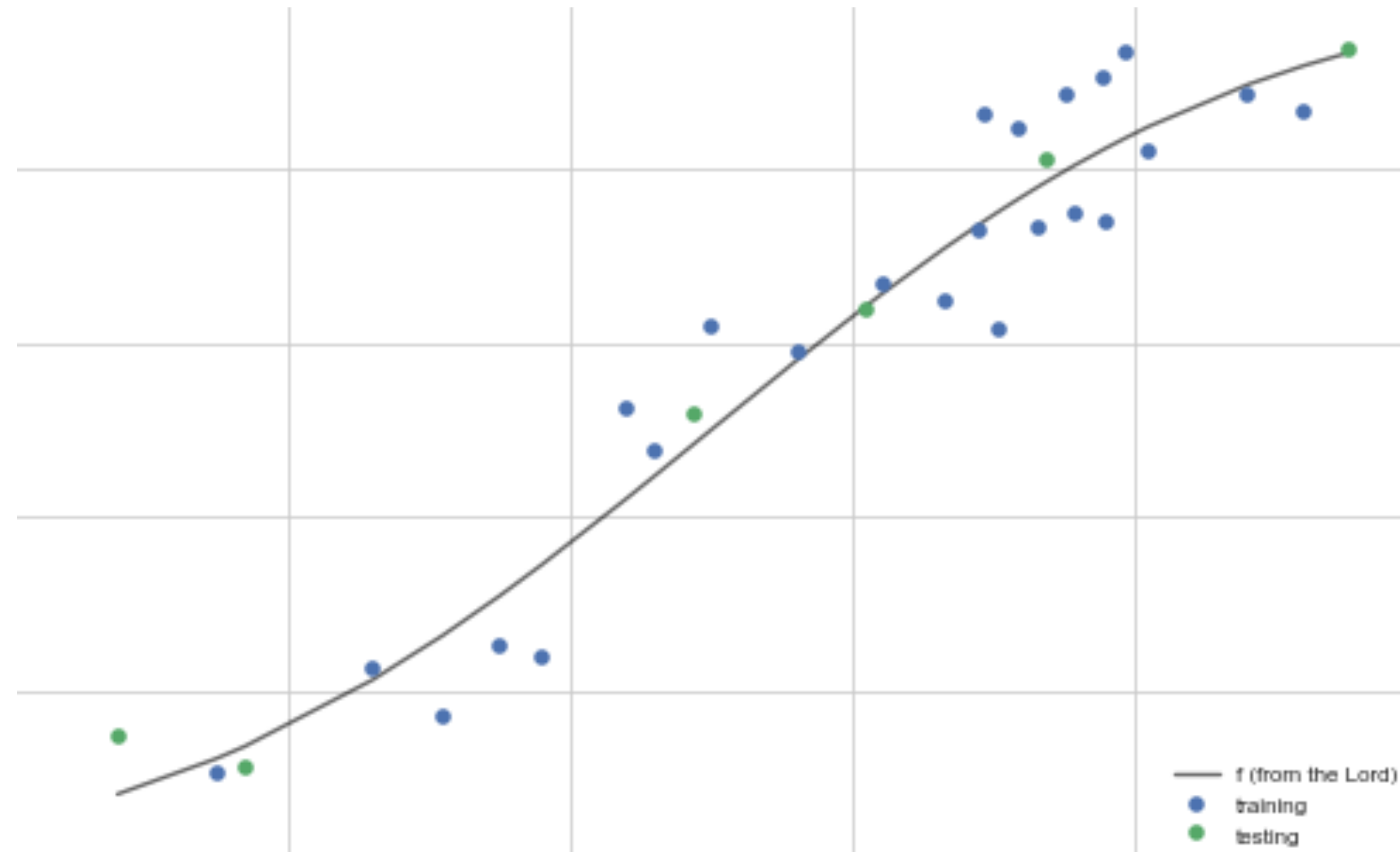
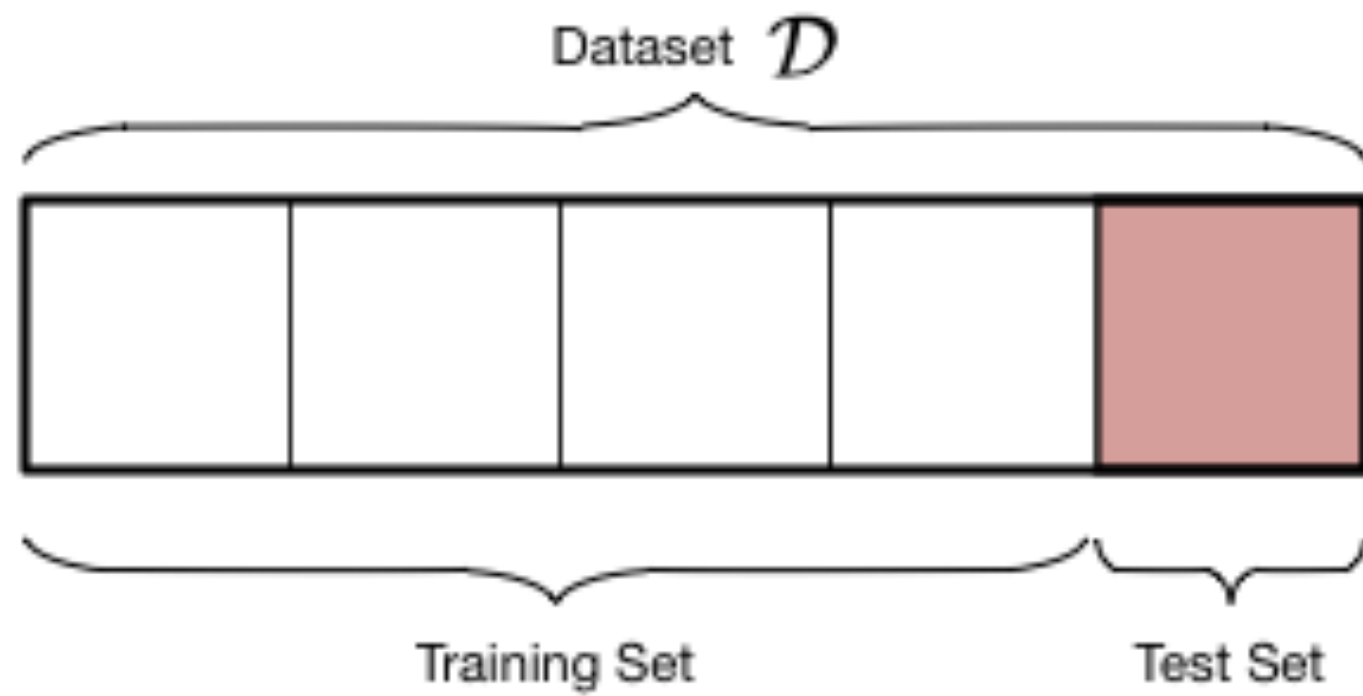
- first term is **variance**, squared error of the various fit g 's from the average g , the hairiness.
- second term is **bias**, how far the average g is from the original f this data came from.
- third term is the **stochastic noise**, minimum error that this model will always have.

DATA SIZE MATTERS: straight line fits to a sine curve

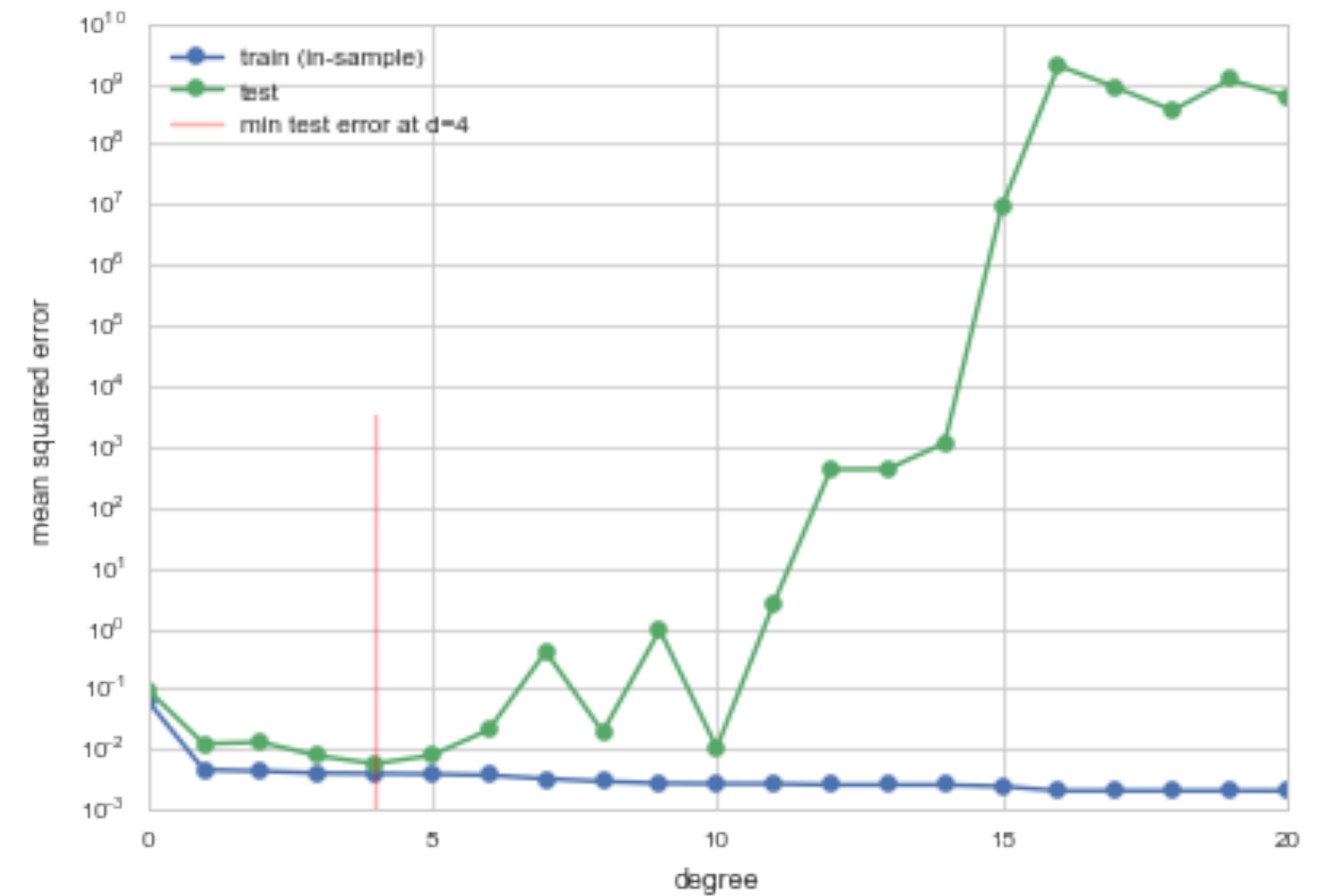
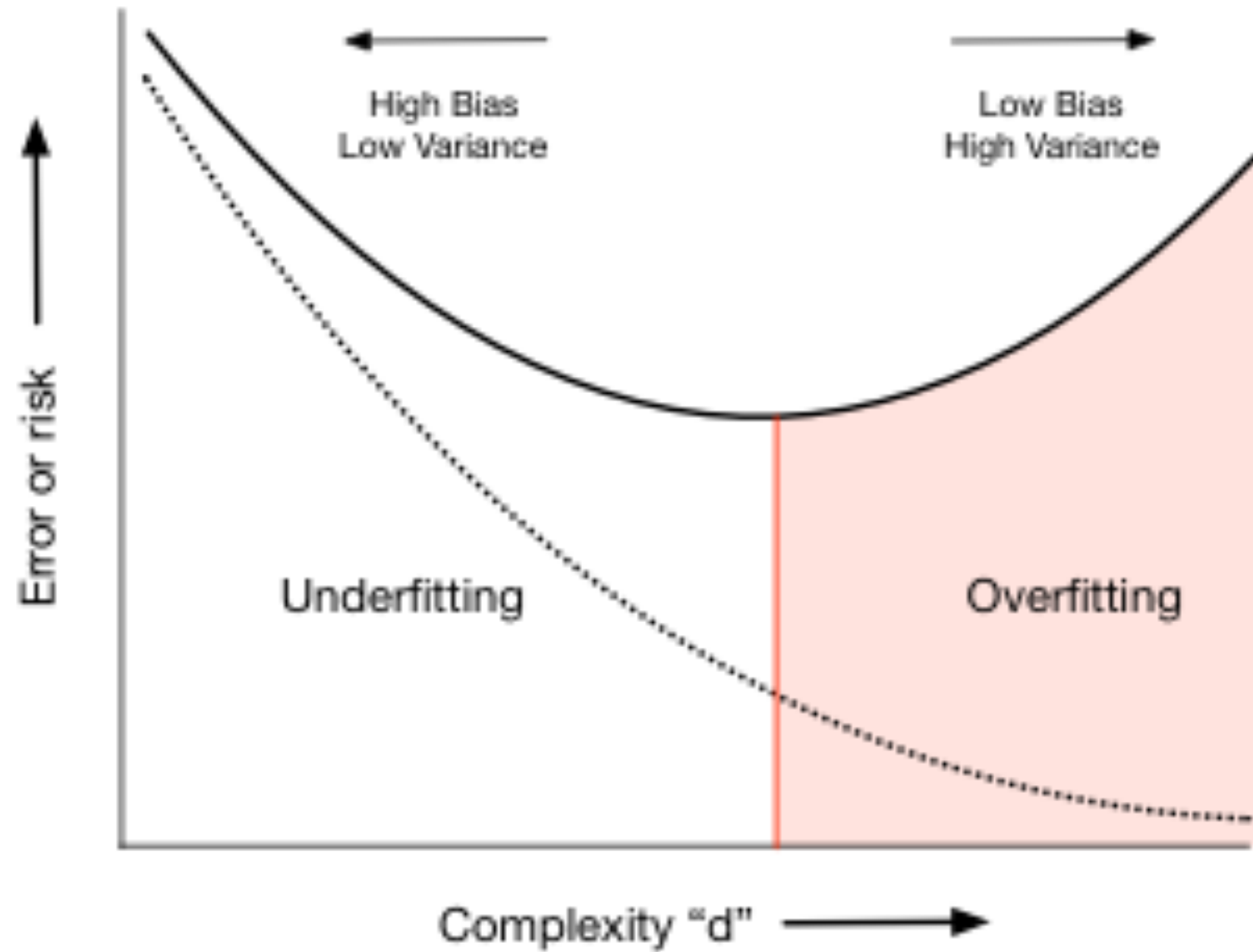


Corollary: Must fit simpler models to less data!

TRAIN AND TEST



BALANCE THE COMPLEXITY



Is this still a test set?

Trouble:

- no discussion on the error bars on our error estimates
- "visually fitting" a value of $d \implies$ contaminated test set.

The moment we **use it in the learning process, it is not a test set.**

Hoeffding's inequality

population fraction μ , sample drawn with replacement, fraction ν :

$$P(|\nu - \mu| > \epsilon) \leq 2e^{-2\epsilon^2 N}$$

For hypothesis h , identify 1 with $h(x_i) \neq f(x_i)$ at sample x_i . Then μ, ν are population/sample error rates. Then,

$$P(|R_{in}(h) - R_{out}(h)| > \epsilon) \leq 2e^{-2\epsilon^2 N}$$

- Hoeffding inequality holds ONCE we have picked a hypothesis h , as we need it to label the 1 and 0s.
- But over the training set we one by one pick all the models in the hypothesis space
- best fit g is among the h in \mathcal{H} , g must be h_1 OR h_2 OR....Say **effectively** M such choices:

$$P(|R_{in}(g) - R_{out}(g)| \geq \epsilon) \leq \sum_{h_i \in \mathcal{H}} P(|R_{in}(h_i) - R_{out}(h_i)| \geq \epsilon) \leq 2 M e^{-2\epsilon^2 N}$$

Hoeffding, rephrased:

Now let $\delta = 2 M e^{-2\epsilon^2 N}$.

Then, **with probability** $1 - \delta$:

$$R_{out} \leq R_{in} + \sqrt{\frac{1}{2N} \ln\left(\frac{2M}{\delta}\right)}$$

For finite effective hypothesis set size M , $R_{out} \sim R_{in}$ as N larger..

Training vs Test

- training error approximates out-of-sample error slowly
- is test set just another sample like the training sample?
- key observation: test set is looking at only one hypothesis because the fitting is already done on the training set. So $M = 1$ for this sample!

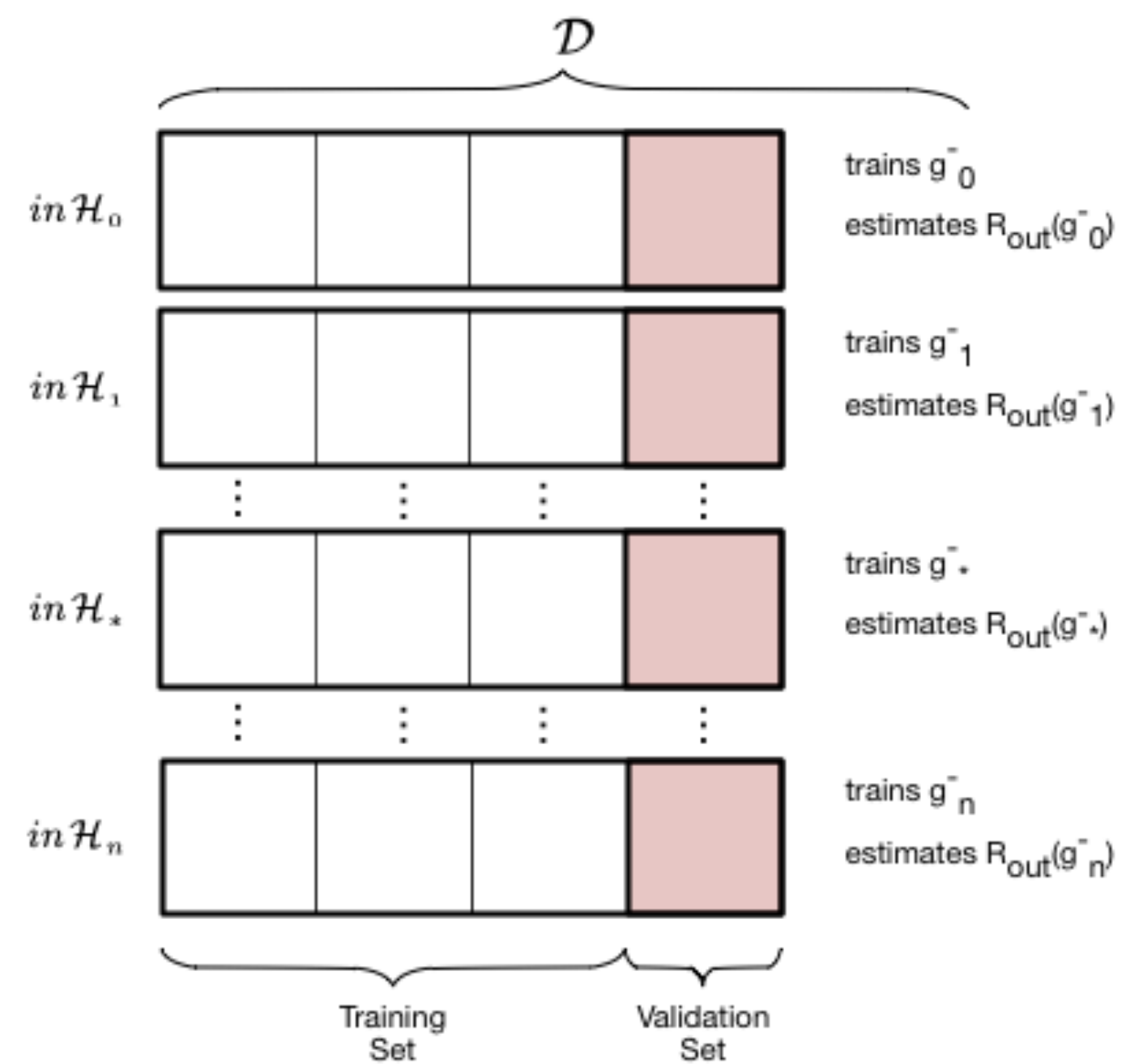
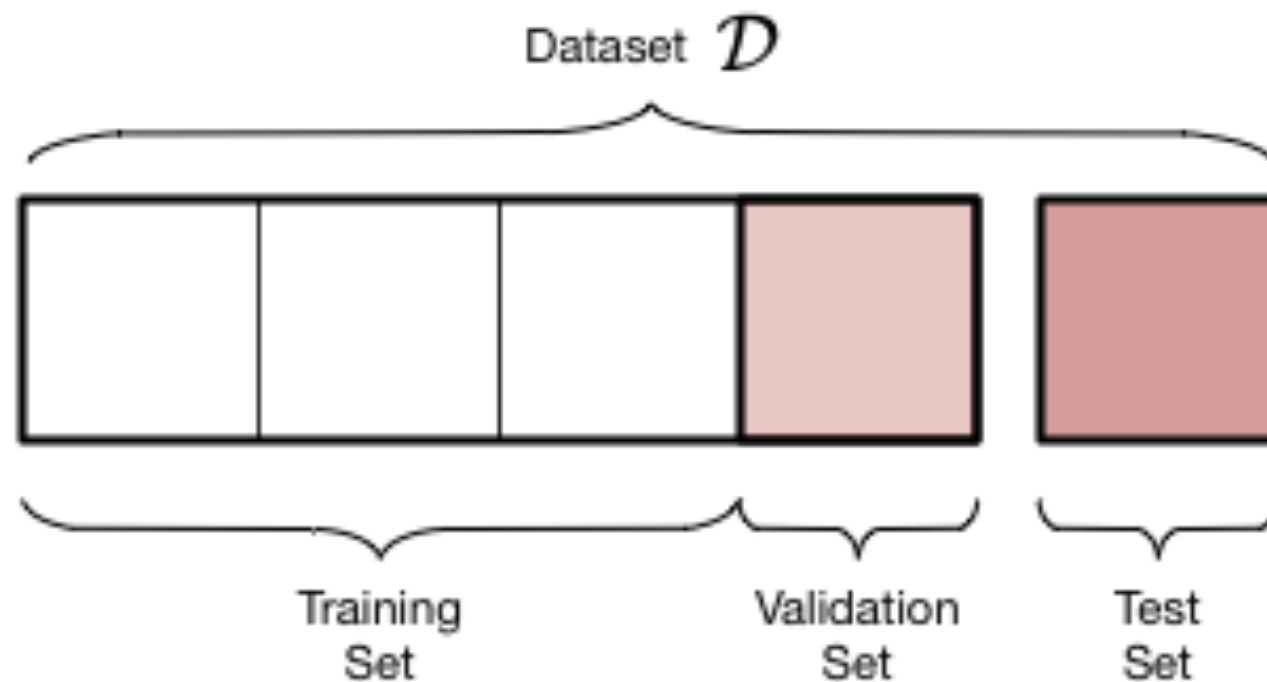
$$R_{out} \leq R_{in} + \sqrt{\frac{1}{2N_{test}} \ln\left(\frac{2}{\delta}\right)}$$

Training vs Test

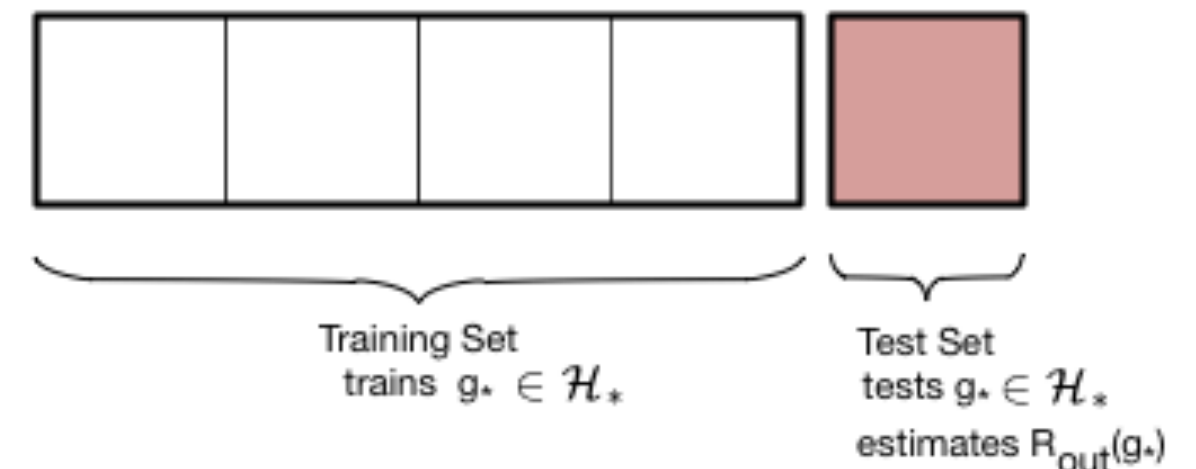
- the test set does not have an optimistic bias like the training set (that's why the larger effective M factor)
- once you start fitting for things like d on the test set, you can't call it a test set any more since we lose tight guarantee.
- test set has a cost of less data in the training set and must thus fit a less complex model.

VALIDATION

- train-test not enough as we *fit* for d on test set and contaminate it
- thus do train-validate-test



pick \mathcal{H}_* with lowest $R_{\text{out}}(g_*^*)$, then retrain in \mathcal{H}_* on entire set



If we don't fit a hyperparameter

- first assume that the validation set is acting like a test set.
- validation risk or error is an unbiased estimate of the out of sample risk.
- Hoeffding bound for a validation set is then identical to that of the test set.

usually we want to fit a hyperparameter

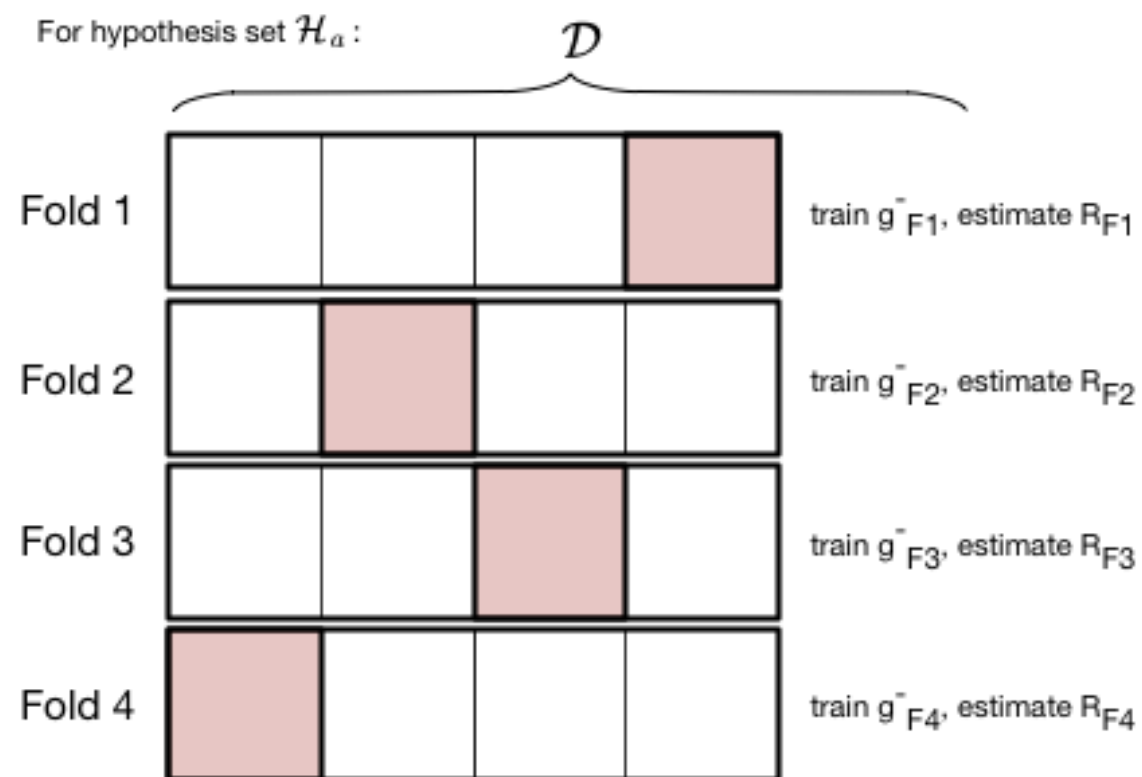
- we **wrongly** already attempted to do on our previous test set.
- choose the d, g^* combination with the lowest validation set risk.
- $R_{val}(g^{-*}, d^*)$ has an optimistic bias since d effectively fit on validation set
- its Hoeffding bound must now take into account the grid-size as the effective size of the hypothesis space.

- this size from hyperparameters is typically a smaller size than that from parameters.

Retrain on entire set!

- finally retrain on the entire train+validation set using the appropriate (g^{-*}, d^{*}) combination.
- works as training for a given hypothesis space with more data typically reduces the risk even further.

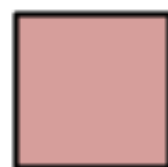
CROSS-VALIDATION



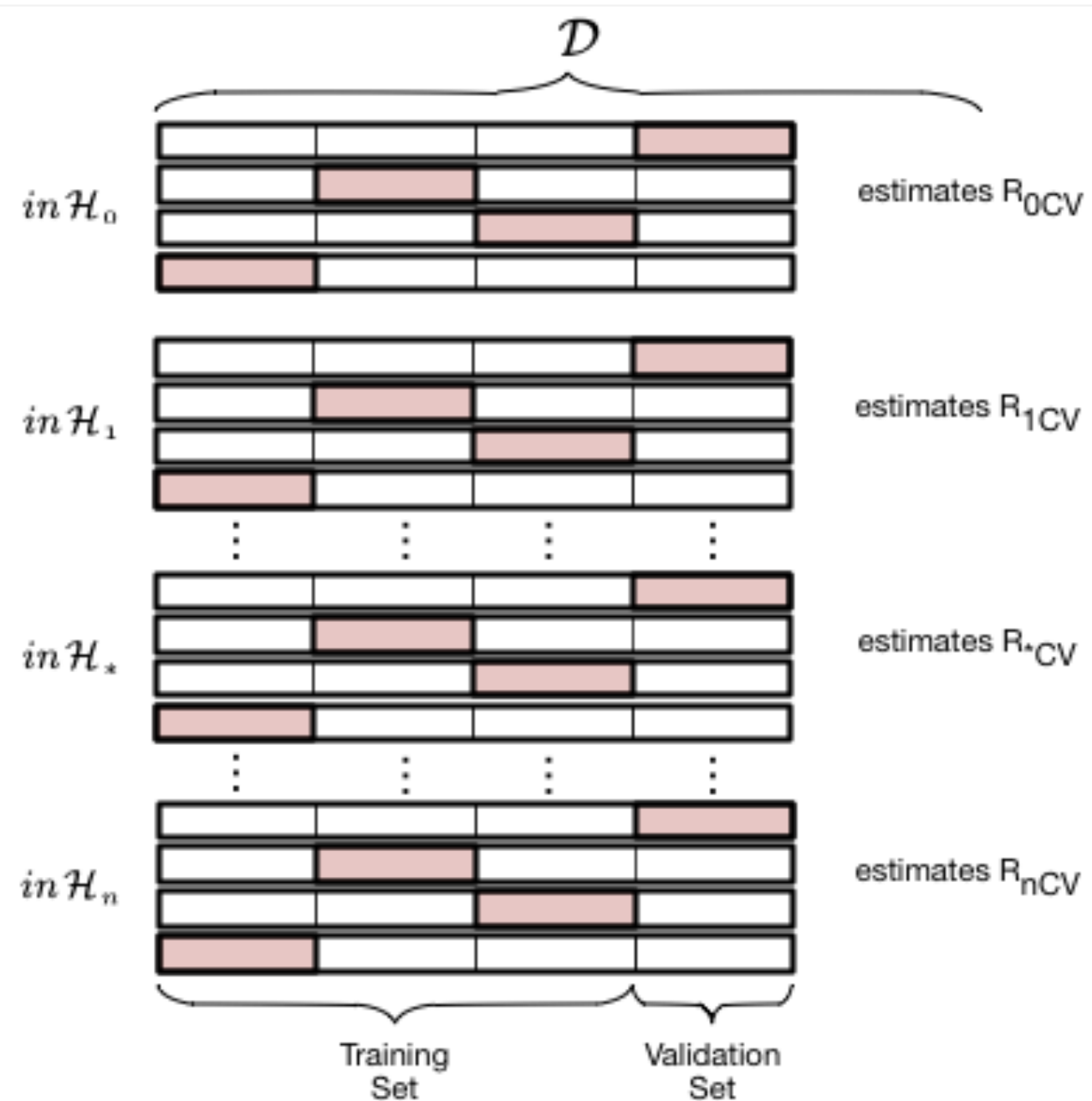
Calculate total error or risk over folds:

$$R_{CV} = \frac{R_{F1} + R_{F2} + R_{F3} + R_{F4}}{4}$$

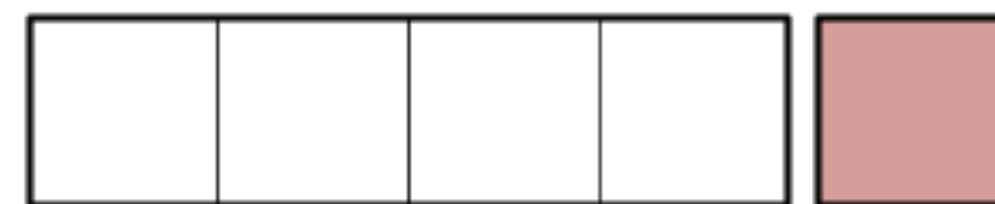
For hypothesis \mathcal{H}_a report R_{CV}



Test Set
left over

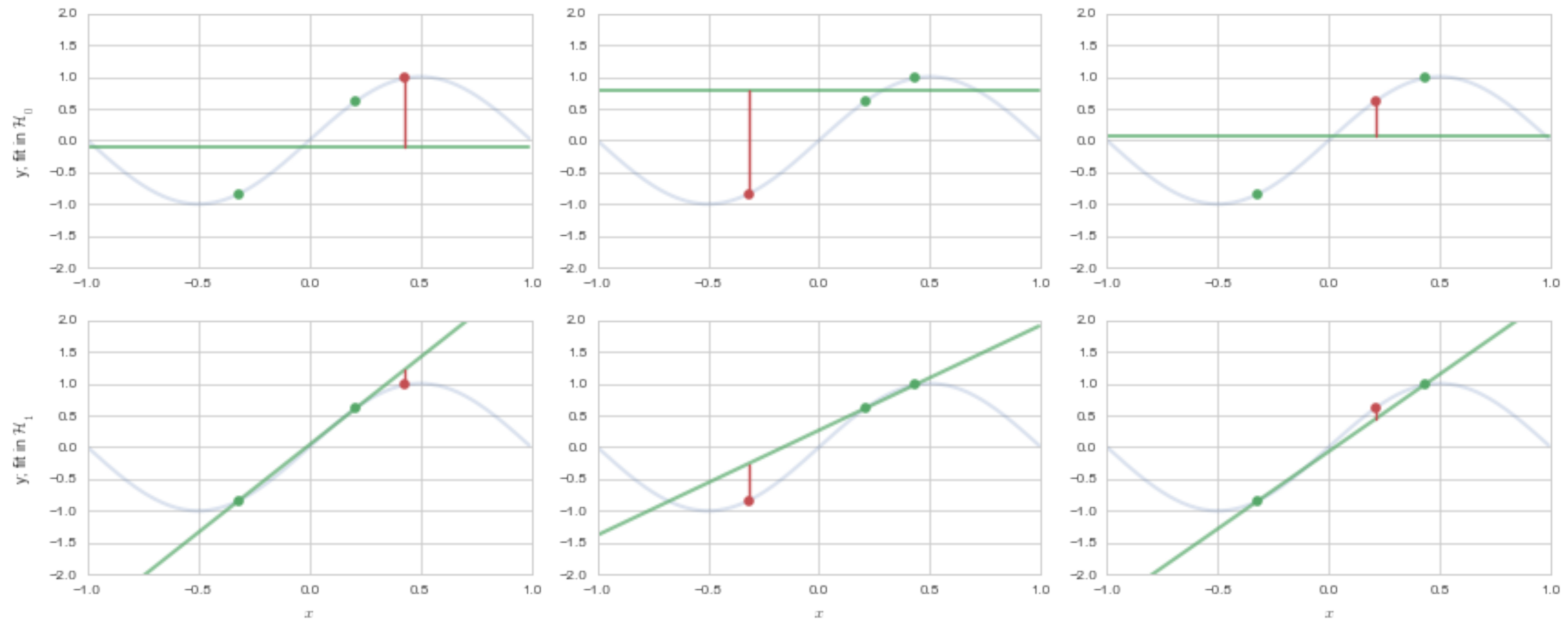


pick \mathcal{H}_* with lowest R_{CV} , then retrain in \mathcal{H}_* on entire set



Training Set
trains $g_* \in \mathcal{H}_*$

Test Set
tests $g_* \in \mathcal{H}_*$
estimates $R_{out}(g_*)$

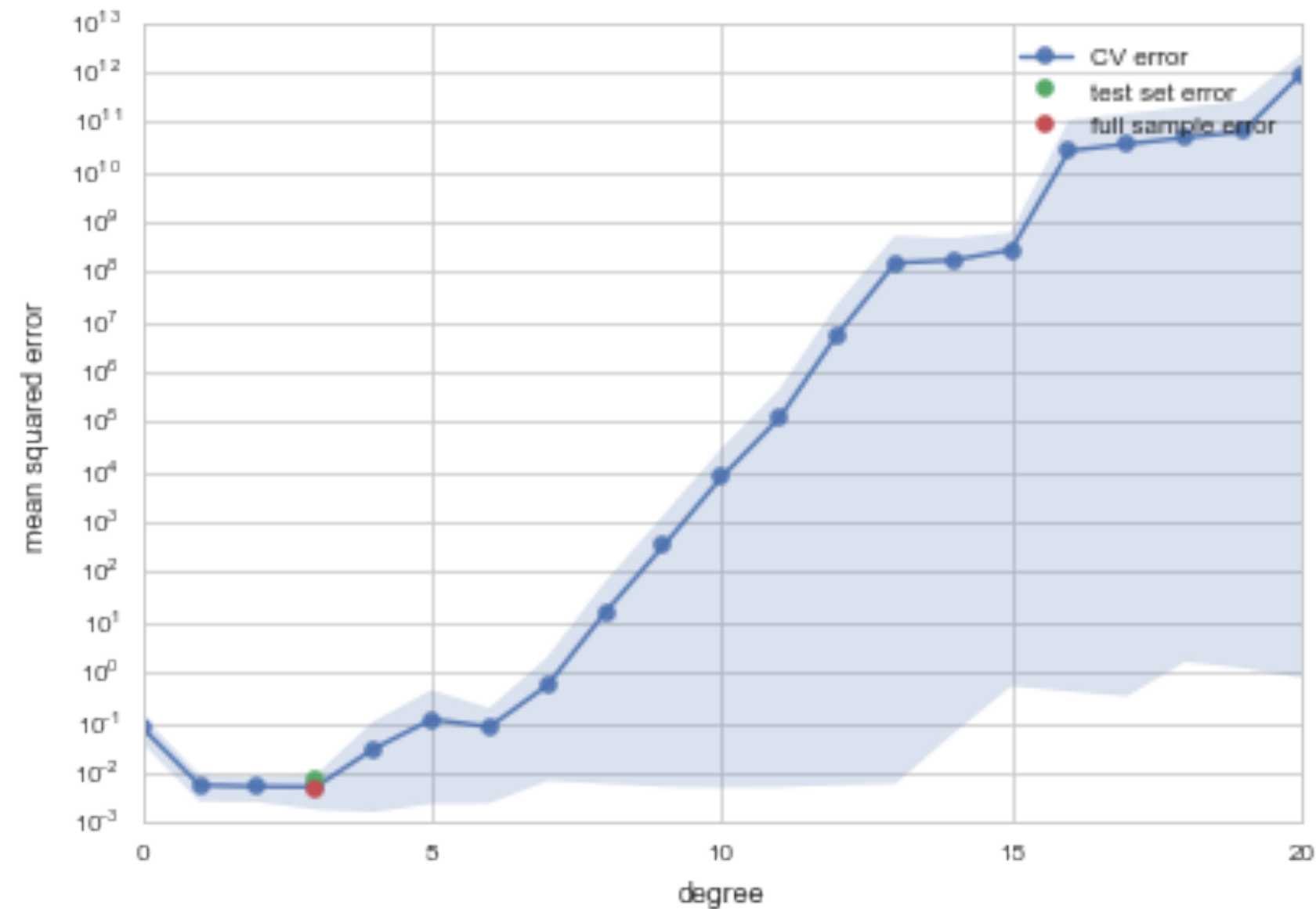


CROSS-VALIDATION

is

- a resampling method
- robust to outlier validation set
- allows for larger training sets
- allows for error estimates

Here we find $d = 3$.



Cross Validation considerations

- validation process as one that estimates R_{out} directly, on the validation set. It's critical use is in the model selection process.
- once you do that you can estimate R_{out} using the test set as usual, but now you have also got the benefit of a robust average and error bars.
- key subtlety: in the risk averaging process, you are actually averaging over different g^- models, with different parameters.

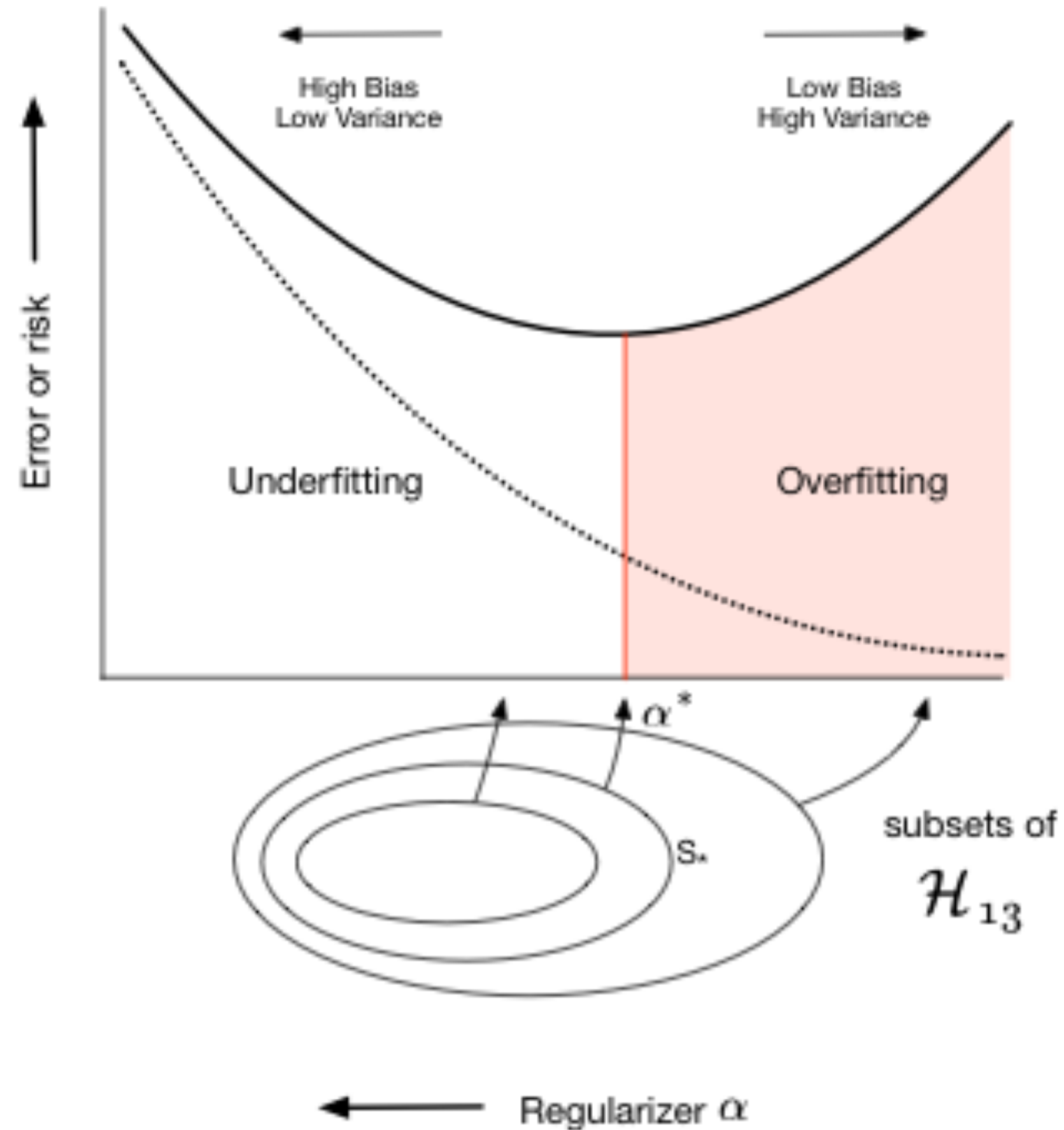
REGULARIZATION

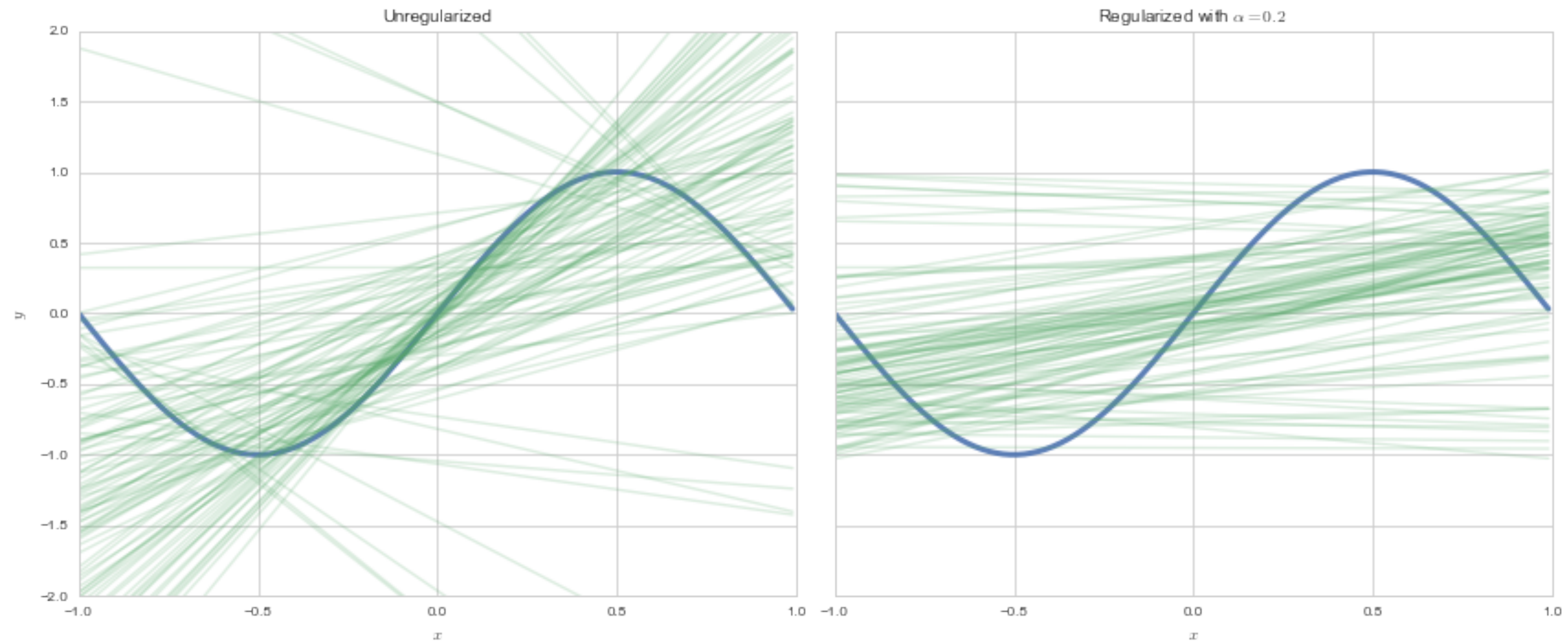
Keep higher a-priori complexity and impose a

complexity penalty

on risk instead, to choose a SUBSET of \mathcal{H}_{big} . We'll make the coefficients small:

$$\sum_{i=0}^j \theta_i^2 < C.$$



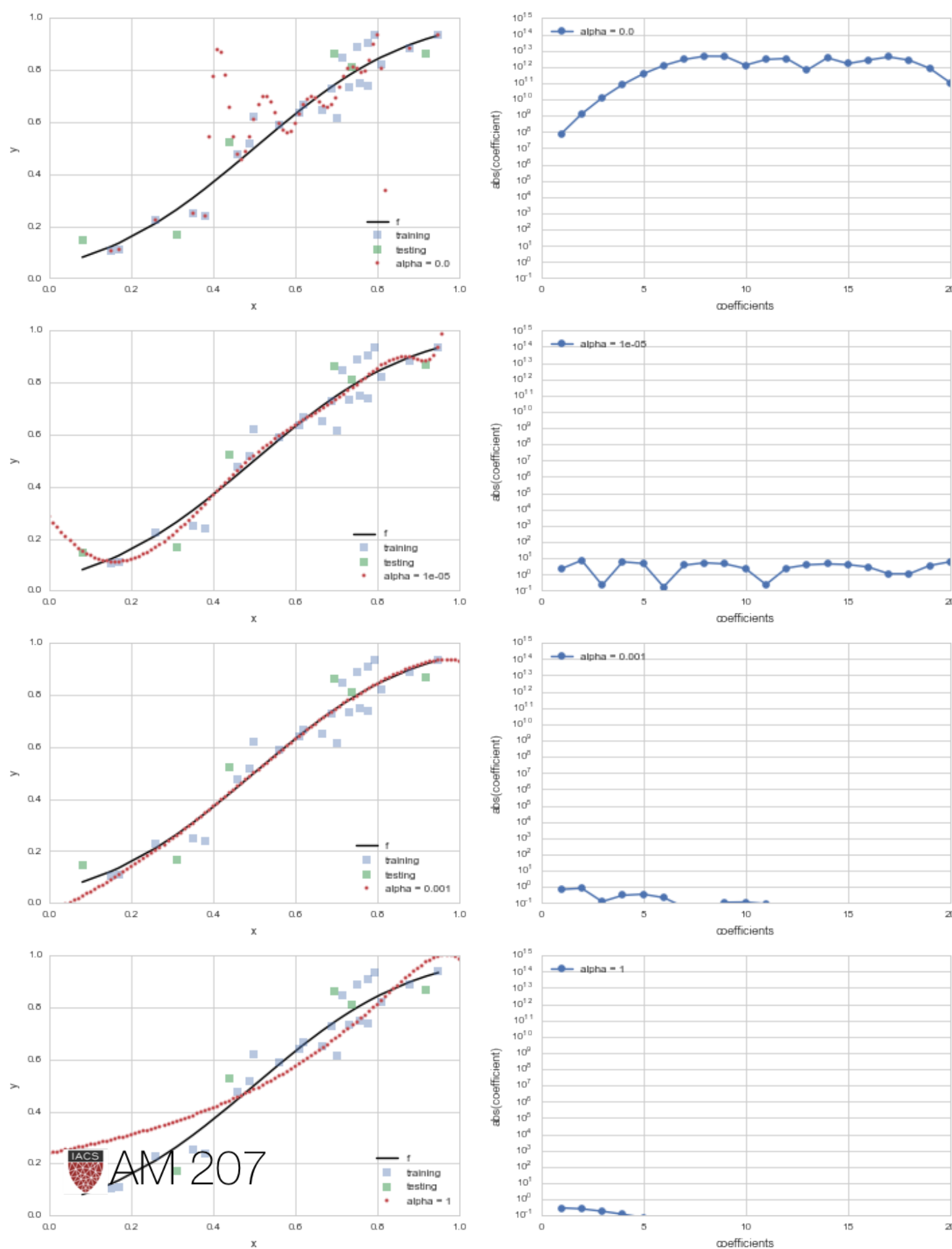


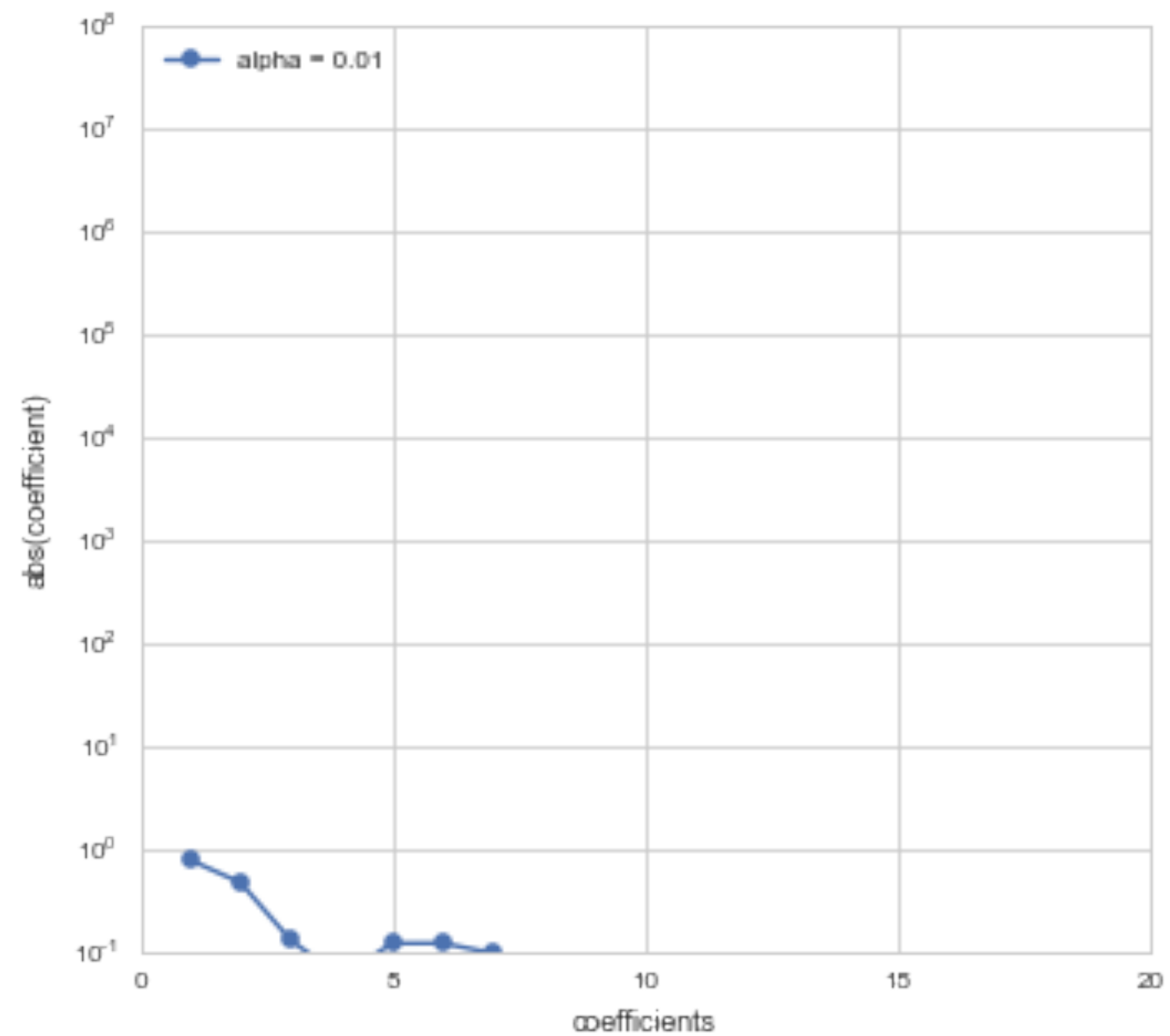
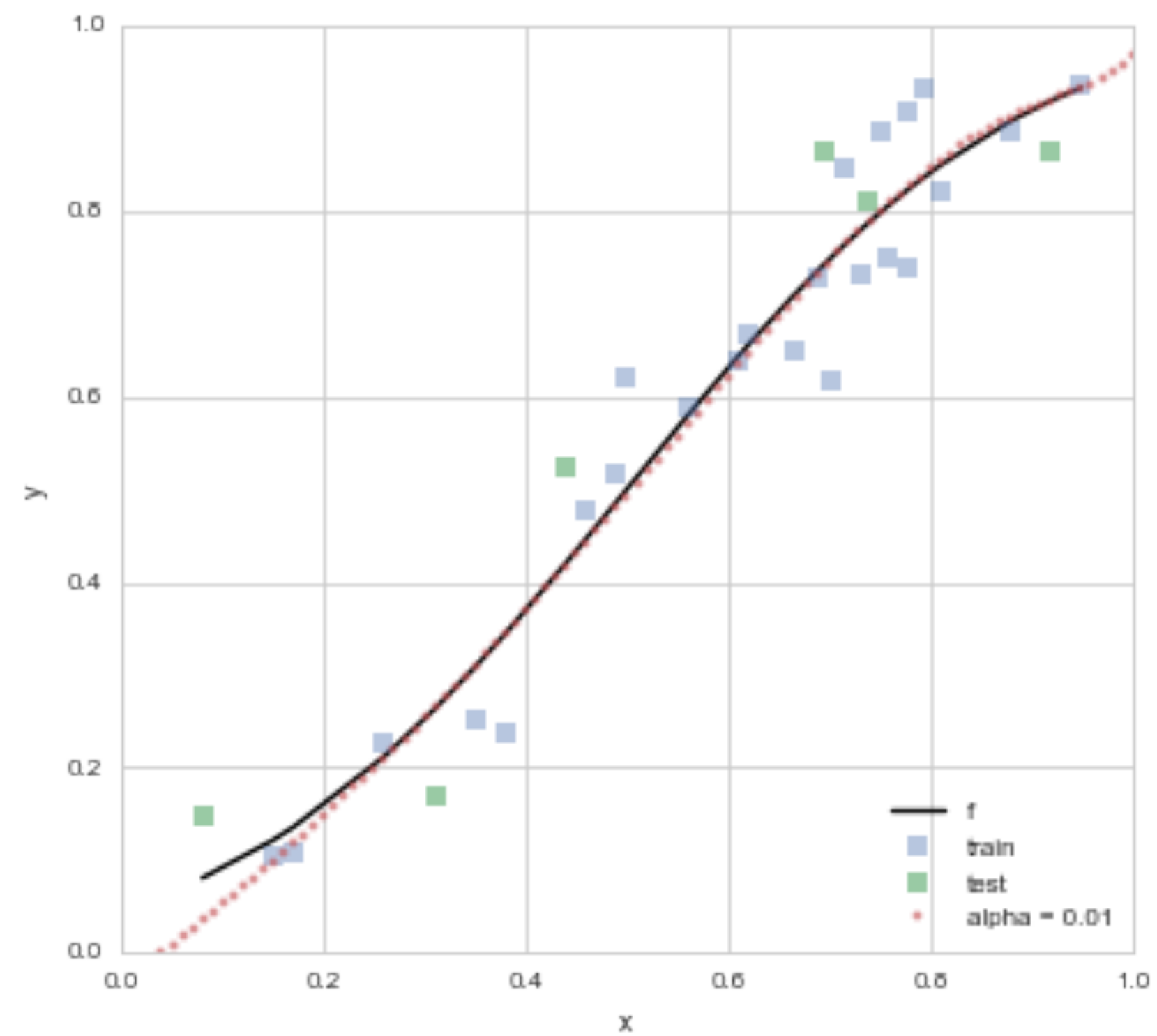
REGULARIZATION

$$\mathcal{R}(h_j) = \sum_{y_i \in \mathcal{D}} (y_i - h_j(x_i))^2 + \alpha \sum_{i=0}^j \theta_i^2.$$

As we increase α , coefficients go towards 0.

Lasso uses $\alpha \sum_{i=0}^j |\theta_i|$, sets coefficients to exactly 0.





Next time

Minimize the risk

- analytically
- using gradient descent
- using stochastic gradient descent