

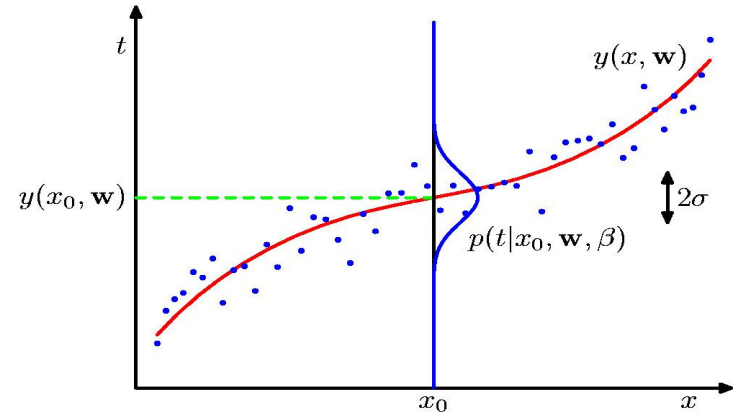
HTF: Ch3, 7
B, Ch3

Linear Regression, Evaluation, Bias-Variance Tradeoff

Thanks to: C Guestrin, T Dietterich, R Parr

Outline

- Linear Regression
 - MLE = Least Squares!
 - Basis functions
- Evaluating Predictors
 - Training set error vs Test set error
 - Cross Validation
- Model Selection
 - Bias-Variance analysis
 - Regularization, Bayesian Model



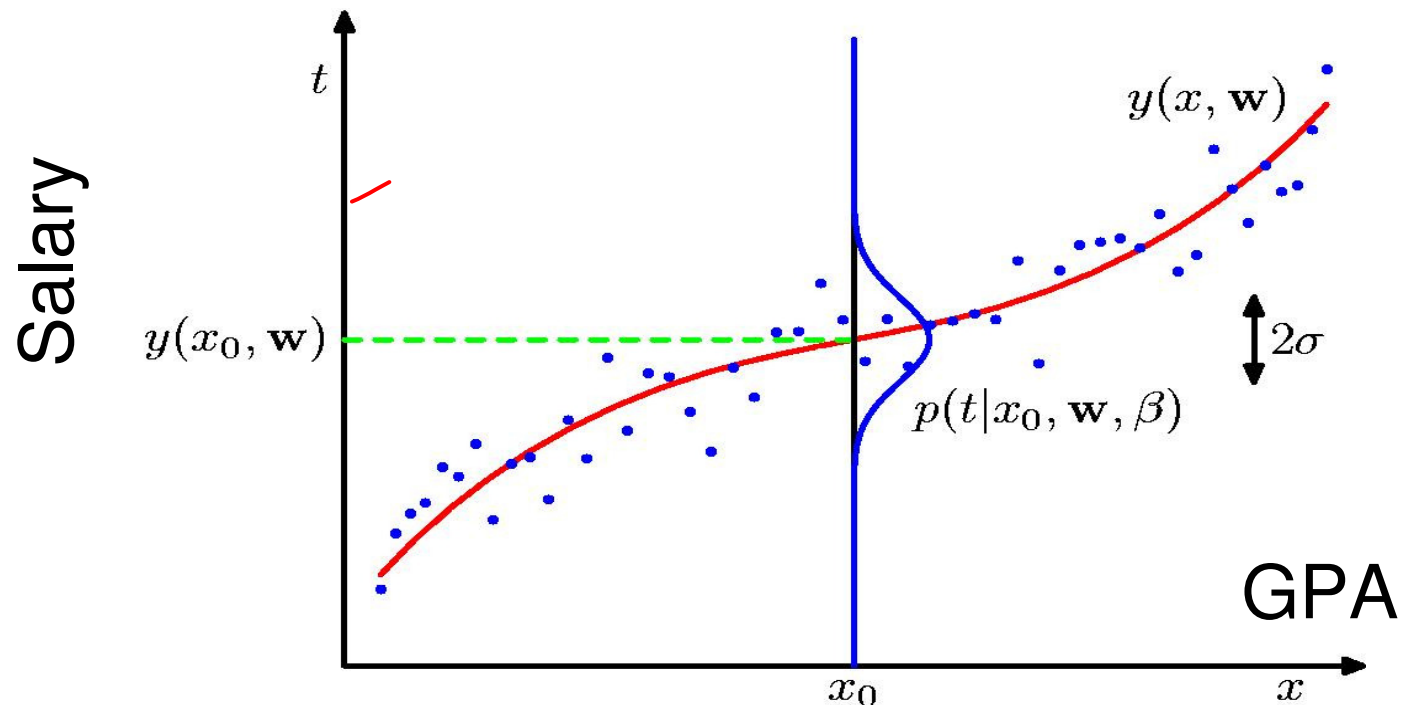


Prediction Problems ...

- Predict housing price from:
 - House size, lot size, rooms, neighborhood, ...
- Predict weight from:
 - Gender, height, ethnicity, ...
- Predict life expectancy increase from:
 - Medication, disease state, ...
- Predict crop yield from:
 - Precipitation, fertilizer, temperature, ...

Prediction of Continuous Variables

- Predict a continuous variable based on set of continuous inputs:
 - Eg, predict salaries from GPA
 - **Regression**

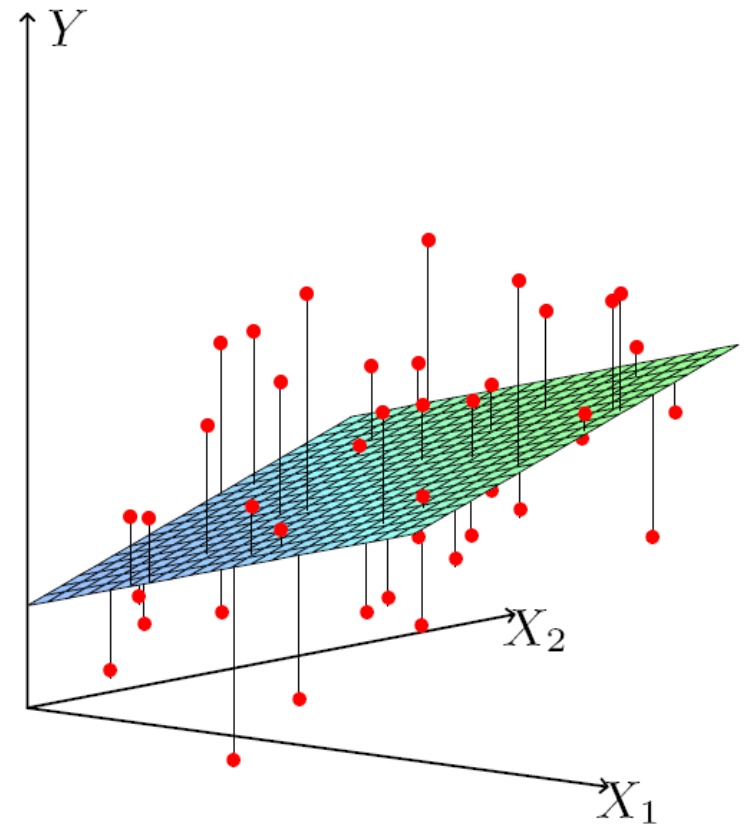


Best LINEAR Fit

- *Finding LINEAR fit*

- *Find $(\beta_0, \beta_1, \dots, \beta_k)$*

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$$



- *Linear least squares fitting with $X \in \mathcal{R}^2$*
- *Seek the linear function of X that minimizes the sum of squared residuals from Y*

The Linear Regression Task

- **Given set of labeled Instances:** $\{ [\mathbf{x}_j, t_j] \}$

GPA, Age, ShoeSize, ... \rightarrow Salary

Eg: [(97, 14, 8); 150]
[(93, 24, 12); 200]
[(88, 20, 9); 45]

- **Learn:** Mapping from \mathbf{x} to $t(\mathbf{x})$

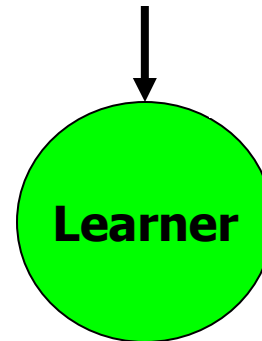
- Direct linear mapping: $t(\mathbf{x}) \approx \beta_0 + \sum_j \beta_j x_j$

- Find coeffs $\beta = (\beta_0, \beta_1, \dots, \beta_k)$

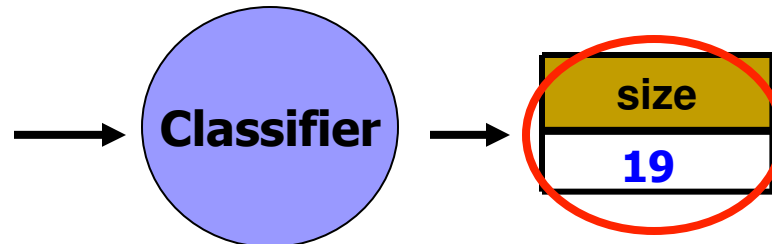
- **Model:** Observed value $t(\mathbf{x}) = \beta_0 + \sum_j \beta_j x_j + \varepsilon$
where $\varepsilon \sim N(0, \sigma^2)$

Training a Regressor

Width	Size	Eyes	...	Light	size
35	95	Y	...	Pale	22
22	110	N	...	Clear	18
:	:			:	:
10	87	N	...	Pale	33



Width	Size	Eyes	...	Light
32	90	N	...	Pale



Best Values of β ?

- **Model:** Observed value is $t(\mathbf{x}) = \sum_j \beta_j x_j + \varepsilon$
where $\varepsilon \sim N(0, \sigma^2)$

$$P(t, \mathbf{x} | \boldsymbol{\beta}, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{[t - \sum_i \beta_i x_i]^2}{2\sigma^2}}$$

- Find MostLikely values of $\boldsymbol{\beta}$... (MLE)

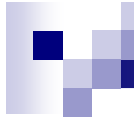
$$\ln P(D | \boldsymbol{\beta}, \sigma)$$

$$= N \ln \left(\frac{1}{\sigma \sqrt{2\pi}} \right) - \frac{1}{2\sigma^2} \left[\sum_j [t^j - \sum_i \beta_i x_i^j]^2 \right]$$

Max Likely Estimate

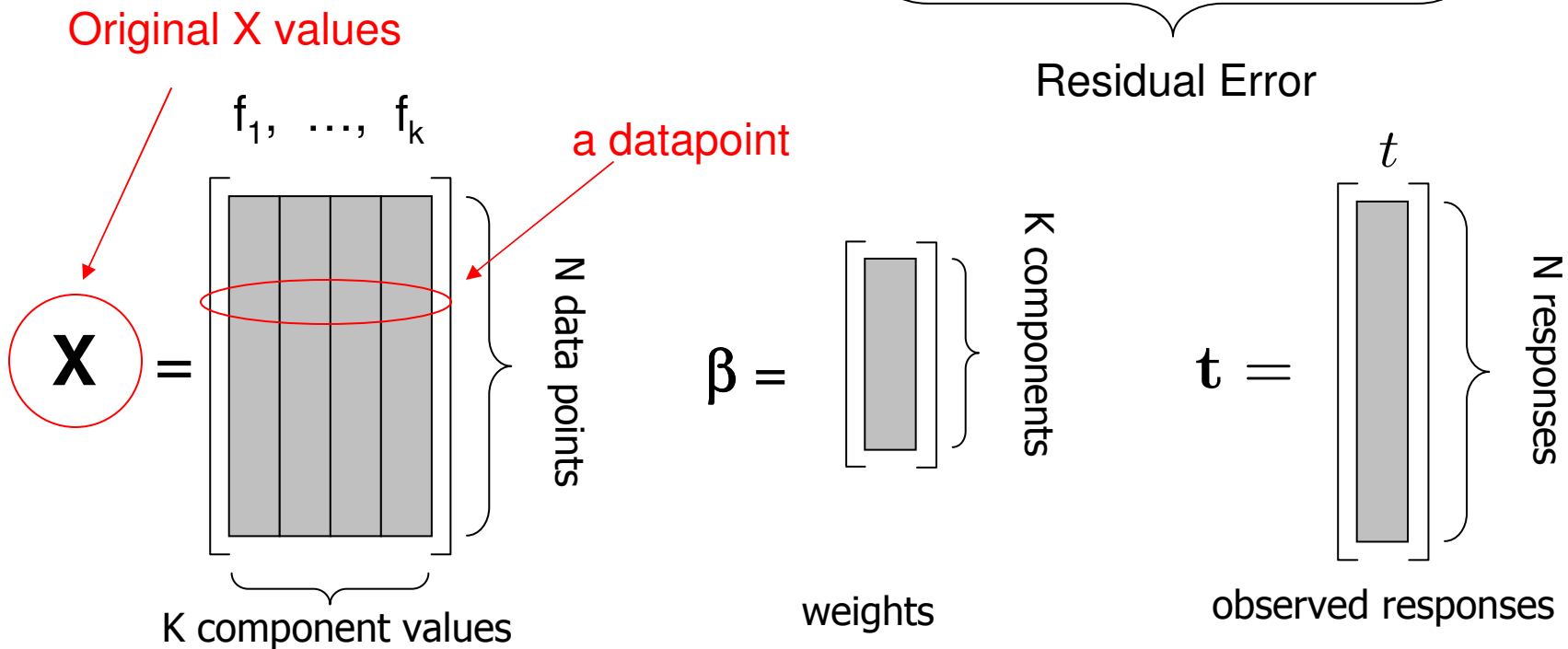
$$\ln P(D | \boldsymbol{\beta}, \sigma) = N \ln \left(\frac{1}{\sigma \sqrt{2\pi}} \right) - \frac{1}{2\sigma^2} \left[\sum_j \left[t^j - \sum_i \beta_i x_i^j \right]^2 \right]$$
$$\arg \max_w \ln P(D | \boldsymbol{\beta}, \sigma) = \arg \min_w \left[\sum_j \left[t^j - \sum_i \beta_i x_i^j \right]^2 \right]$$

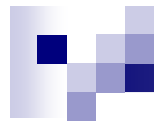
- **Least-squares Linear Regression is MLE for Gaussians !!!**



Regression in Matrix Notation

$$\beta^* = \arg \min_{\beta} \sum_{j=1}^N \left[t^j - \sum_{i=1}^K \beta_i x_i^j \right]^2$$
$$= \arg \min_{\beta} \underbrace{(\mathbf{X}\beta - \mathbf{t})^T (\mathbf{X}\beta - \mathbf{t})}_{\text{Residual Error}}$$



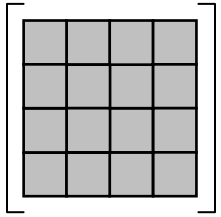
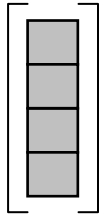


Regression solution = simple matrix operations

$$\boldsymbol{\beta}^* = \arg \min_{\boldsymbol{\beta}} (\mathbf{X}\boldsymbol{\beta} - \mathbf{t})^T (\mathbf{X}\boldsymbol{\beta} - \mathbf{t})$$

Setting derivative to 0 yields:

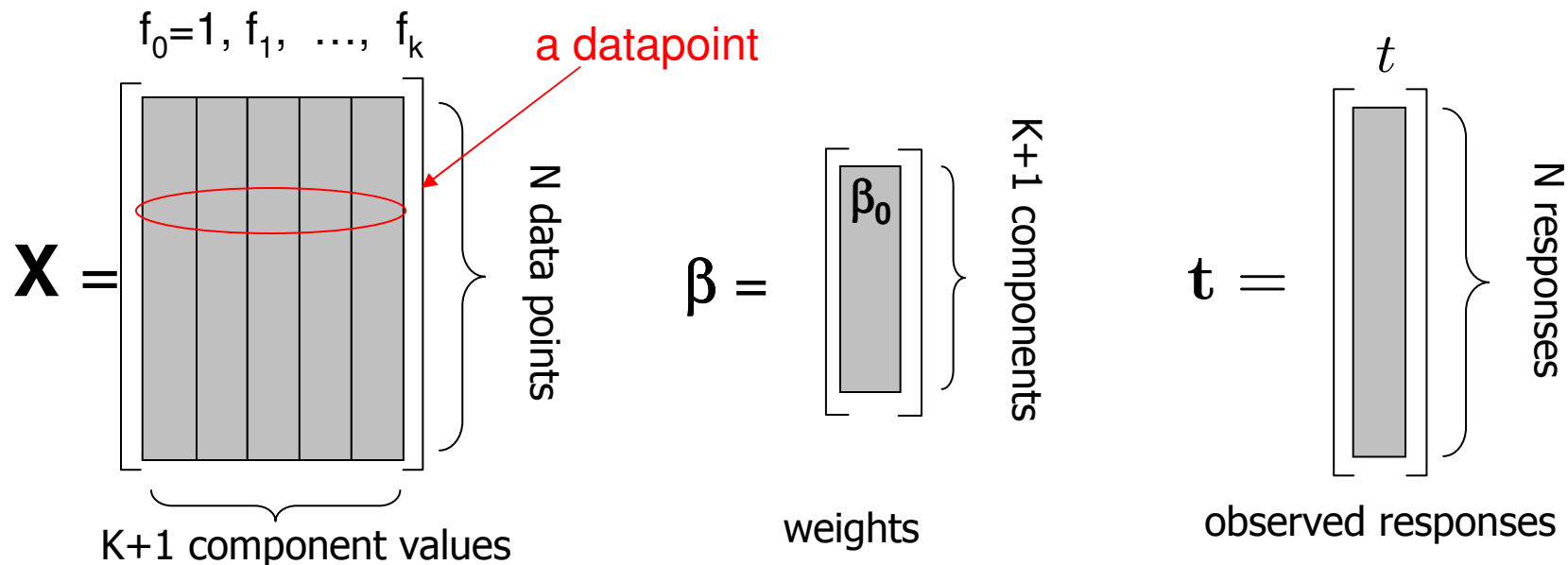
$$\text{Solution: } \boldsymbol{\beta}^* = \underbrace{(\mathbf{X}^T \mathbf{X})^{-1}}_{\mathbf{A}^{-1}} \underbrace{\mathbf{X}^T \mathbf{t}}_{\mathbf{b}} = \mathbf{A}^{-1} \mathbf{b}$$

where $\mathbf{A} = \mathbf{X}^T \mathbf{X} =$  $\mathbf{b} = \mathbf{X}^T \mathbf{y} =$ 

$\underbrace{\hspace{10em}}_{\text{K} \times \text{K matrix for k components}$ $\underbrace{\hspace{10em}}_{\text{K} \times 1 \text{ vector}$

Dealing with Offset

- Actually want $k+1$ values $(\beta_0, \beta_1, \dots, \beta_k)$
$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$$
- So view each k -tuple \mathbf{x} as $k+1$ tuple $[1, \mathbf{x}]$

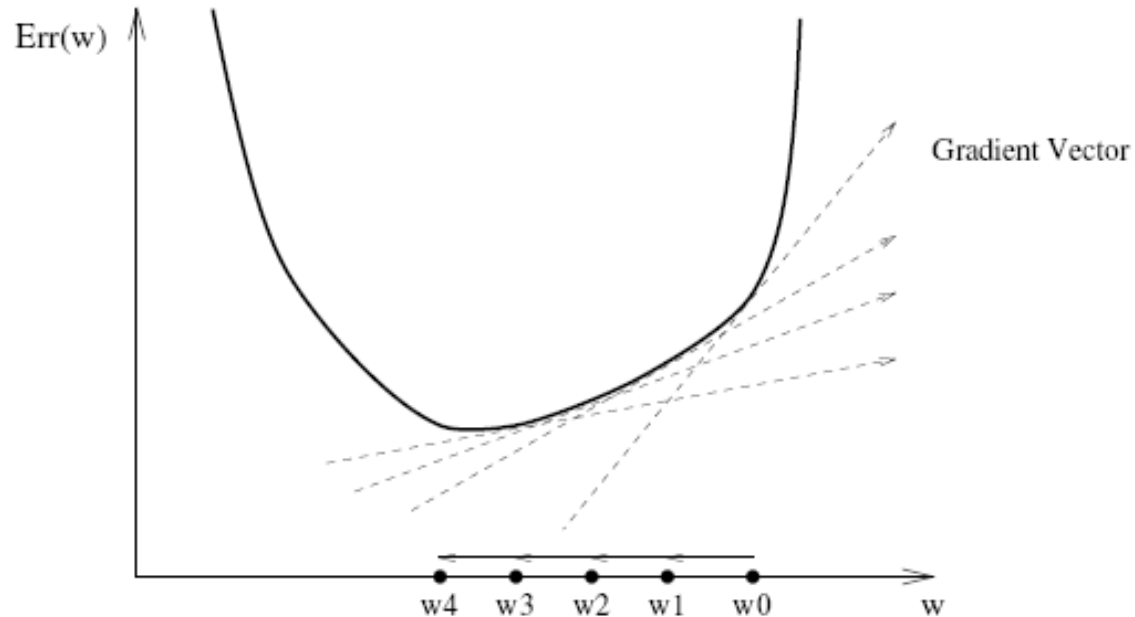


Another Approach: Gradient Descent

Skip 2009

- But matrix solution is $O(m^4d)$
 - $m = \text{\#training instances}$; $d = \text{\#features}$
 - matrix inversion...
- Goal: Find w_i 's that minimize squared error
 - $\text{err}(w) = (\sum_i [t^{(i)} - w x^{(i)}]^2) / m$
- Why not use Gradient Descent!
 - aka Delta Rule, Adaline Rule, Widrow-Ho Rule, LMS Rule, Classical Conditioning

Local Search via Gradient Descent



Start w / (random) weight vector w^0 .
Repeat until converged \vee bored

Compute Gradient

$$\nabla \text{err}(w^t) = \left(\frac{\partial \text{err}(w^t)}{\partial w_0}, \frac{\partial \text{err}(w^t)}{\partial w_1}, \dots, \frac{\partial \text{err}(w^t)}{\partial w_n} \right)$$

Let $w^{t+1} = w^t + \eta \nabla \text{err}(w^t)$

If CONVERGED: Return(w^t)



Gradient-Descent vs Matrix-Inversion

Pro: Gradient Descent advantages

- ≈Biologically plausible
- Each iteration costs only $O(mn)$
- If uses $< m$ iterations, faster than Matrix Inversion!
- More easily parallelizable

Con: Gradient Descent disadvantages

- It's moronic... essentially a slow way to build $X^T X$ matrix, then solve a set of linear equations
- If n is small, it's especially outrageous. If n is large then direct matrix inversion method can be problematic but not impossible if you want to be efficient.
- Need to choose a good learning rate -- how?
- Matrix inversion takes predictable time. You can't be sure when gradient descent will stop.

Computing the Gradient

$$\text{err}(w) = (\sum_i [t^{(i)} - w x^{(i)}]^2) / m$$

$$\frac{\partial \text{err}(w)}{\partial w_j} = \frac{\partial}{\partial w_j} \left(\frac{1}{m} \sum_{i=1}^m \text{err}_i(w) \right) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial w_j} \text{err}_i(w)$$

$$\begin{aligned} \frac{\partial \text{err}_i(w)}{\partial w_j} &= \frac{\partial}{\partial w_j} \left[\left(\sum_{i=1}^m w_j x_{ij} \right) - t_i \right]^2 \\ &= 2 \cdot \left[\left(\sum_{i=1}^m w_j x_{ij} \right) - t_i \right] \cdot \frac{\partial}{\partial w_j} \left[\left(\sum_{i=1}^m w_j x_{ij} \right) - t_i \right] \\ &= 2 \cdot \left[\left(\sum_{i=1}^m w_j x_{ij} \right) - t_i \right] \cdot x_{ij} \end{aligned}$$

Then descend a distance η along gradient $\left[\frac{\partial \text{err}(w)}{\partial w_0}, \frac{\partial \text{err}(w)}{\partial w_1}, \dots, \frac{\partial \text{err}(w)}{\partial w_n} \right]$

Gradient Descent Algorithm

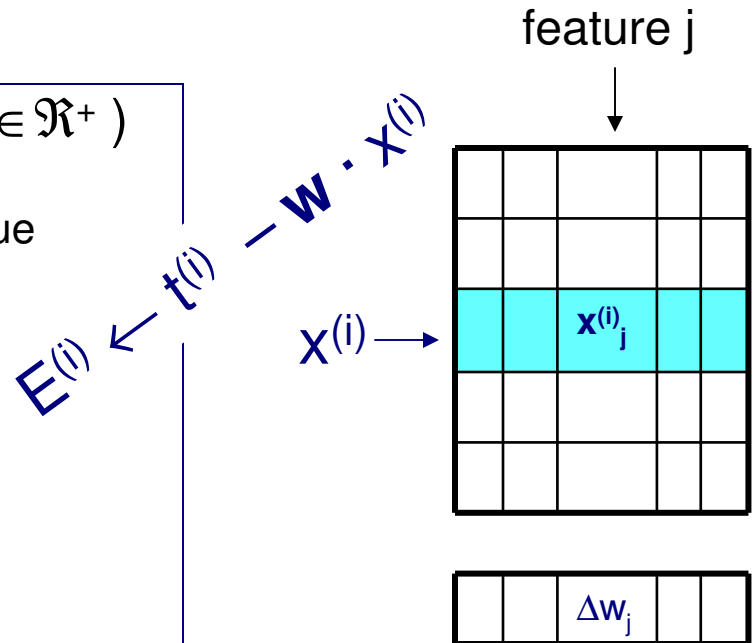
Gradient-Descent(S : training examples; $\eta \in \mathbb{R}^+$)

% $S = \{ [x^{(i)}, t^{(i)}] \}, \dots$

% x = vector of input values; t is target output value

% η is learning rate (eg, 0.05)

- Initialize each w_j to small random value
 - Typically $\in [-0.05, +0.05]$
- Until termination condition is met, do
 - Initialize each $\Delta w_j \leftarrow 0$
 - For each $[x^{(i)}, t^{(i)}] \in S$, do
 - Set $E^{(i)} \leftarrow t^{(i)} - w \cdot x^{(i)}$
 - For each j , do
 - $\Delta w_j \leftarrow \Delta w_j + r^{(i)} x^{(i)}_j$
 - For each j do
 - $w_j \leftarrow w_j + \eta \Delta w_j$
- Return w



0. New \mathbf{w}

1. For each row i , compute

a. $\Delta \mathbf{w} = 0$

b. $\mathbf{E}^{(i)} = \mathbf{t}^{(i)} - \mathbf{w} \cdot \mathbf{x}^{(i)}$

c. $\Delta \mathbf{w} += \mathbf{E}^{(i)} \mathbf{x}^{(i)}$

[... $\Delta w_j += E^{(i)} x_j^{(i)}$...]

2. Increment $\mathbf{w} += \eta \Delta \mathbf{w}$

feature j



$\mathbf{x}^{(i)}$ →

		$\mathbf{x}^{(i)}_j$		

$\mathbf{E}^{(i)}$

$\Delta \mathbf{w}$ →

		Δw_j		
--	--	--------------	--	--

Batch Gradient Descent:

```
Do until satisfied
1. Compute gradient  $\nabla \text{err}_S[\mathbf{w}]$ 
2.  $\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla \text{err}_S[\mathbf{w}]$ 
```

On-Line Gradient Descent:

```
Do until satisfied
  For each training example  $\langle \mathbf{x}, y \rangle$  in  $S$ 
    1. Compute gradient  $\nabla \text{err}_{(\mathbf{x}, y)}[\mathbf{w}]$ 
    2.  $\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla \text{err}_{(\mathbf{x}, y)}[\mathbf{w}]$ 
```

$$\begin{aligned} \text{err}_{(\mathbf{x}, y)}[\mathbf{w}] &\equiv (y - \mathbf{w} \cdot \mathbf{x})^2 \\ \text{err}_S[\mathbf{w}] &\equiv \sum_{(\mathbf{x}, y) \in S} \text{err}_{(\mathbf{x}, y)}[\mathbf{w}] = \sum_{(\mathbf{x}, y) \in S} (y - \mathbf{w} \cdot \mathbf{x})^2 \end{aligned}$$

On-Line Gradient Descent can approximate
Batch Gradient Descent arbitrarily closely
if η small enough



Comments on On-Line Mode

- Stochastic gradient descent may avoid local minima as “noisier”
 - aka “Stochastic Gradient Descent”,
“Robbins-Munro algorithm”
- Some versions store all training examples; make repeated passes through them until convergence

Learning Rates and Convergence

- Learning rate $\eta \equiv$ “step size”
- Convergence whenever...
 - $\lim_{t \rightarrow \infty} \eta_t = 0$
 - $\sum \eta_t = \infty$
 - $\sum \eta_t^2 < \infty$
- \exists sophisticated alg's
(Newton's method; Line Search; ...)
that choose step size automatically, converge faster.
- \exists only one “basin” for linear threshold units
 \Rightarrow local minimum is global minimum!
- Good starting point \Rightarrow algorithm converges faster



Results wrt Gradient Descent

Gradient descent (Delta training rule)

- guaranteed to converge to hypothesis with minimum squared error (eventually!)

if

- Sufficiently small learning rate η
- . . . even when training data
 - contains noise
 - not separable!

What about other features?

■ Data:

- Not linear!!
- Perhaps

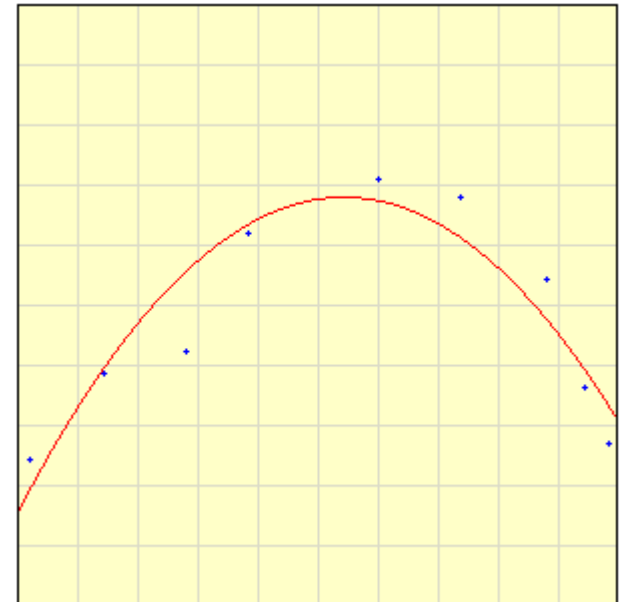
$$f(x) = \alpha_2 x^2 + \alpha_1 x + \alpha_0$$

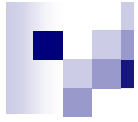
■ How to fit ???

t	x	x²	x	1
18	5	25	5	1
11	-2	4	-2	1
38	7	49	7	1
6	3	9	3	1

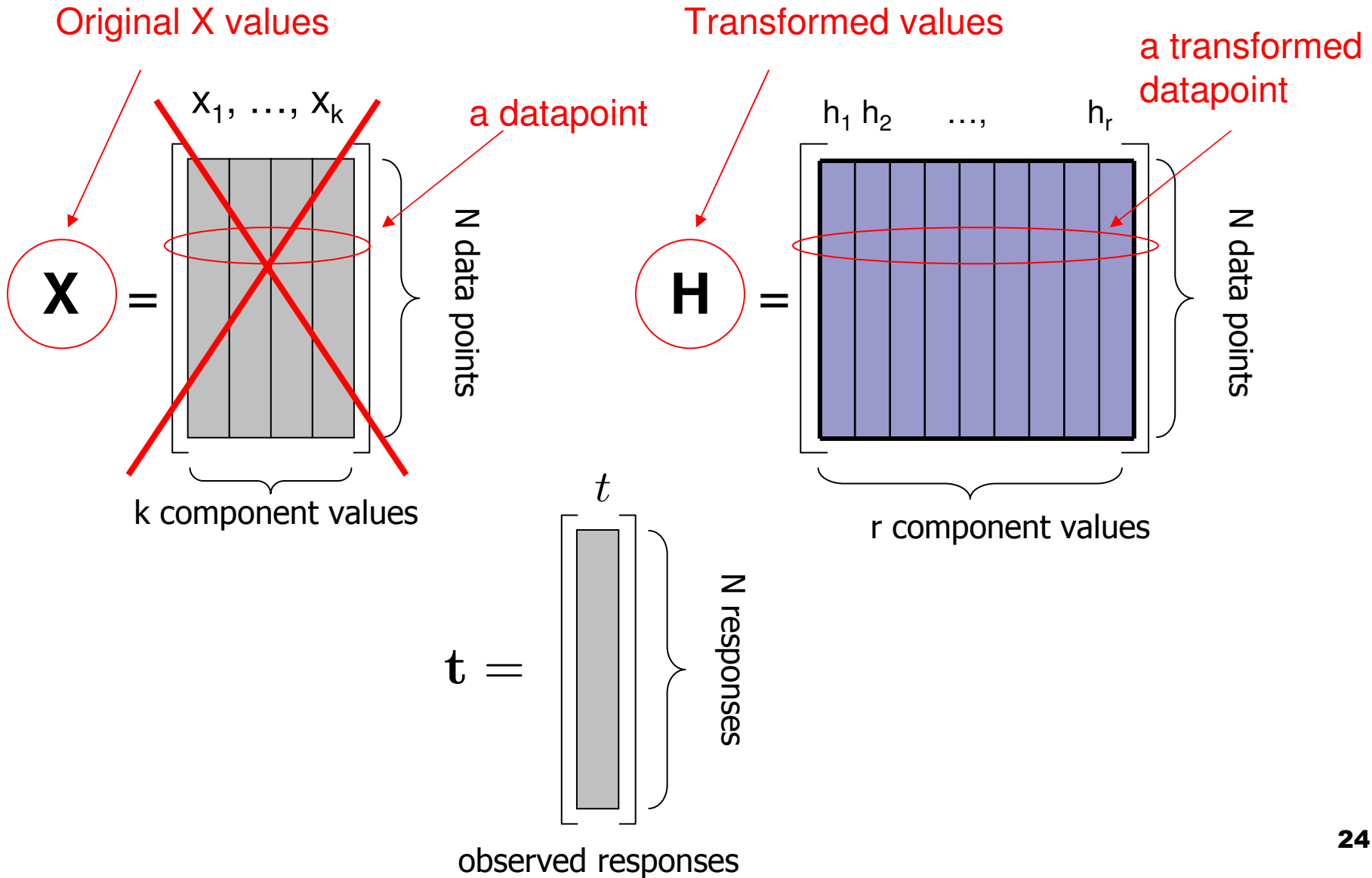
$$\begin{aligned}\alpha_2 &= 1 \\ \alpha_1 &= -2 \\ \alpha_0 &= 3\end{aligned}$$

$$\begin{aligned}\alpha_2 &= -0.179 \\ \alpha_1 &= 1.938 \\ \alpha_0 &= 1.543\end{aligned}$$





General Approach



General Linear Regression Task

- **Given set of labeled Instances:** $\{ [\mathbf{x}_j, t_j] \}$
- **Learn:** Mapping from \mathbf{x} to $t(\mathbf{x})$
 - Can use **BASIS** functions: $H = \{ h_1(\mathbf{x}), \dots, h_r(\mathbf{x}) \}$
 - Eg: $x_i^2, x_i^3, (x_1 x_3), x_i \sin(x_i), \dots$
 - (Basis) linear mapping: $t(\mathbf{x}) \approx \sum_j \beta_j h_j(\mathbf{x})$
 - Find coeffs $\beta = (\beta_1, \dots, \beta_r)$
- **Model:** Observed value $t^*(\mathbf{x}) = \sum_j \beta_j h_j(\mathbf{x}) + \varepsilon$
where $\varepsilon \sim N(0, \sigma^2)$

Model is LINEAR in these bases... even if bases are NOT linear

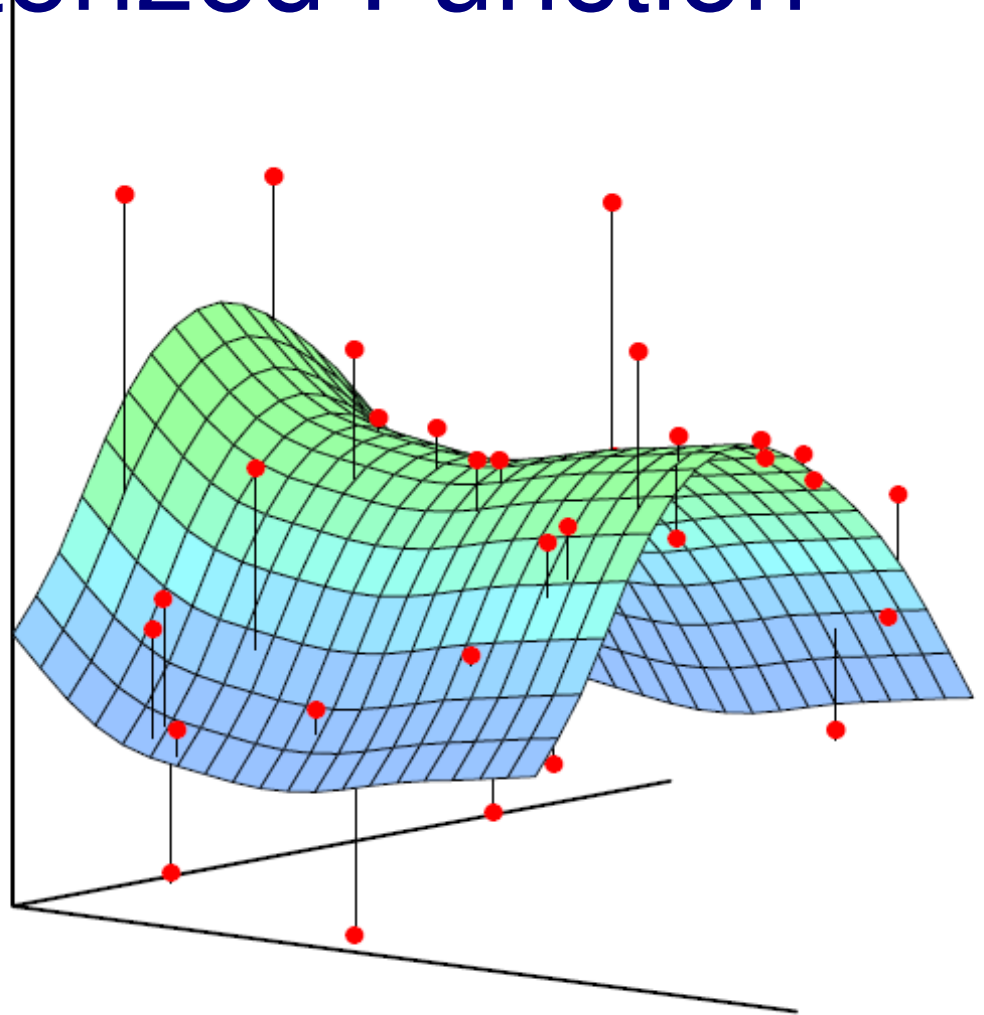


Features/Basis Functions

- Polynomials
 - $1, x, x^2, x^3, x^4, ..$
- Indicators
- Gaussian densities
- Step functions or sigmoids
- Sinusoids (Fourier basis)
- Wavelets
- Anything you can imagine...

Fitting Parameterized Function

- *Least squares fitting of a function of two inputs*
- *Find parameters of $f_{\theta}(x)$ that minimize the sum-of-squared vertical errors*



What if $\mathbf{t}^{(i)}$ is a vector?

- Nothing changes!
- Scalar prediction:

$$\textit{Solution: } \boldsymbol{\beta}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}$$

- Vector prediction:

$$\textit{Solution: } \boldsymbol{\beta}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{T}$$

Target MATRIX

Applications Corner 1

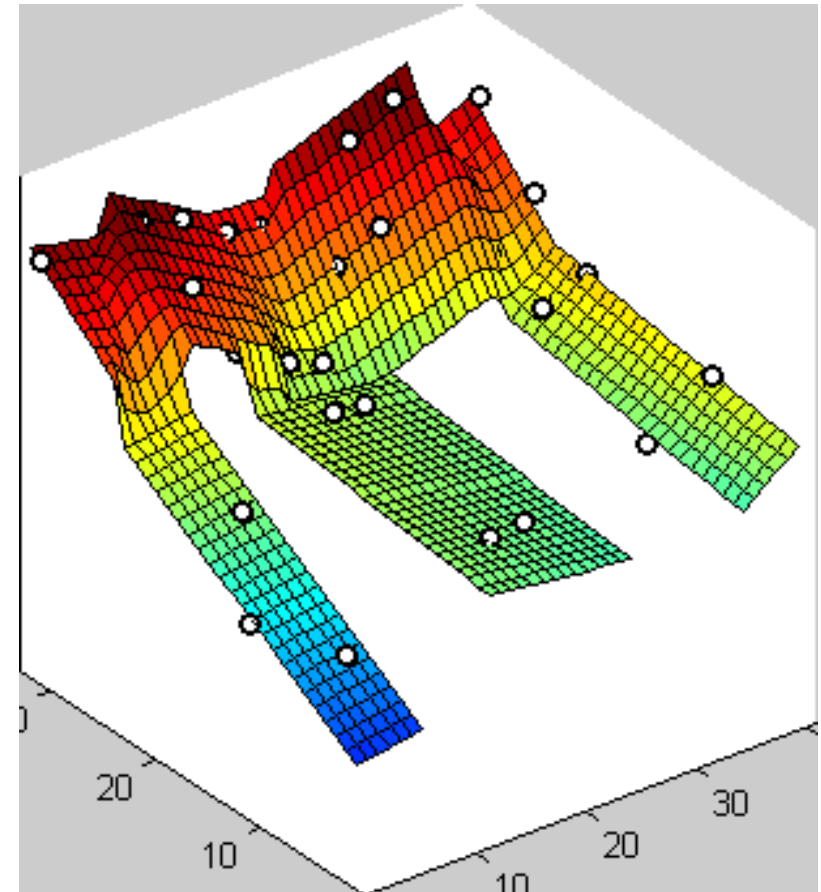
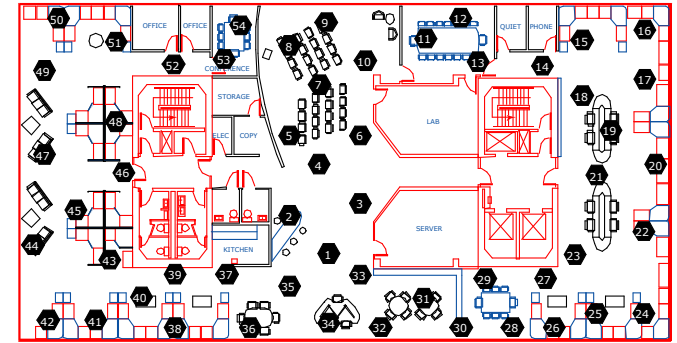
- Predict stock value over time from
 - past values
 - other relevant vars
 - e.g., weather, demands, etc.



Applications Corner 2



- Measure temperatures at some locations
- Predict temperatures throughout the environment



[Guestrin et al. '04]



Applications Corner 3

- Predict when a sensor will fail
 - based several variables
 - age, chemical exposure, number of hours used,...

Outline

- Linear Regression

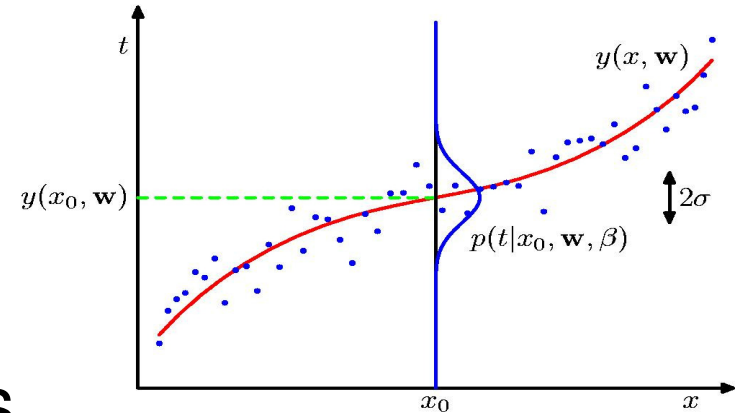
- MLE = Least Squares.
- Basis functions

- Evaluating Predictors

- Training set error vs Test set error
- Cross Validation

- Model Selection

- Bias-Variance analysis
- Regularization, Bayesian Model



Training Set Error

- Choose a loss function
 - eg, squared error (L_2) for regression
- Given a labeled dataset S , learn optimal predictor \mathbf{w}_S

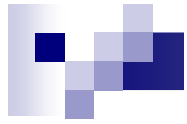
$$\mathbf{w}_S = \mathbf{w}^*(S) = \arg \min_{\mathbf{w}} \sum_{(\mathbf{x}, t) \in S} \left(t - \sum_i w_i h_i(\mathbf{x}) \right)^2$$

S = "Training data"

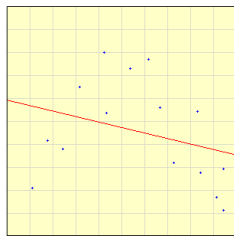
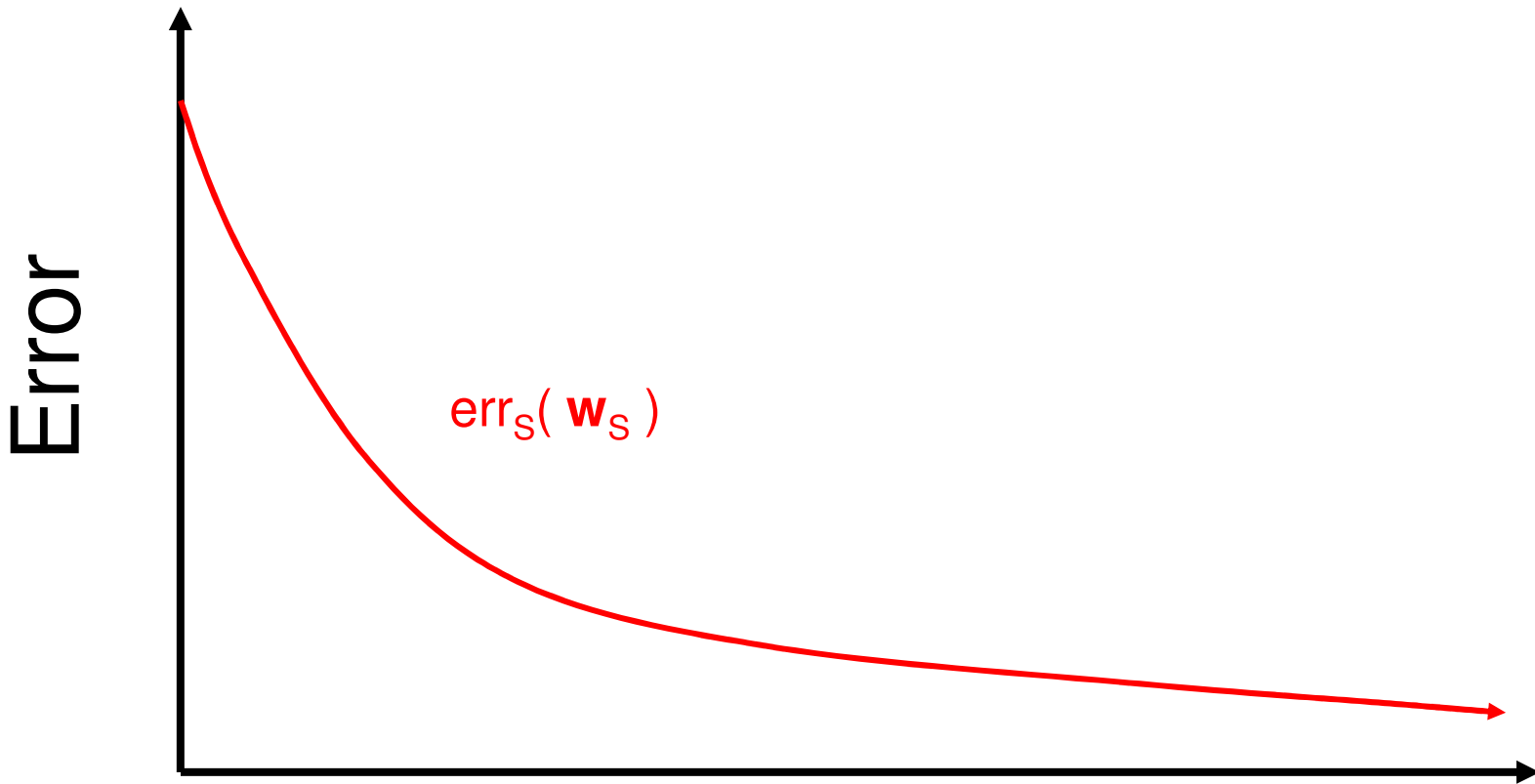
- compute empirical error (of any \mathbf{w})

$$\text{err}_S(\mathbf{w}) = \frac{1}{|S|} \sum_{(\mathbf{x}, t) \in S} \left(t - \sum_i w_i h_i(\mathbf{x}) \right)^2$$

- Training set error: $\text{err}_S(\mathbf{w}_S)$

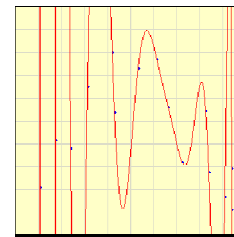


Training Set Error as a function of Model Complexity



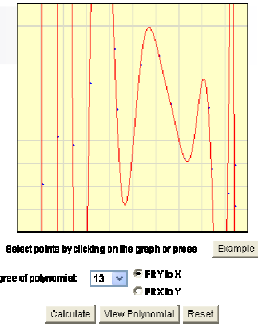
Select points by clicking on the graph or press [Example](#)
Degree of polynomial: FIT Y to X
 FIT X to Y

“Model Complexity”
... eg, #basis functions;
degree of poly, ...



Select points by clicking on the graph or press [Example](#)
Degree of polynomial: FIT Y to X
 FIT X to Y

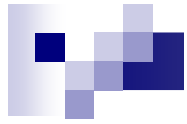
True Prediction Error



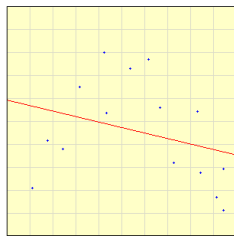
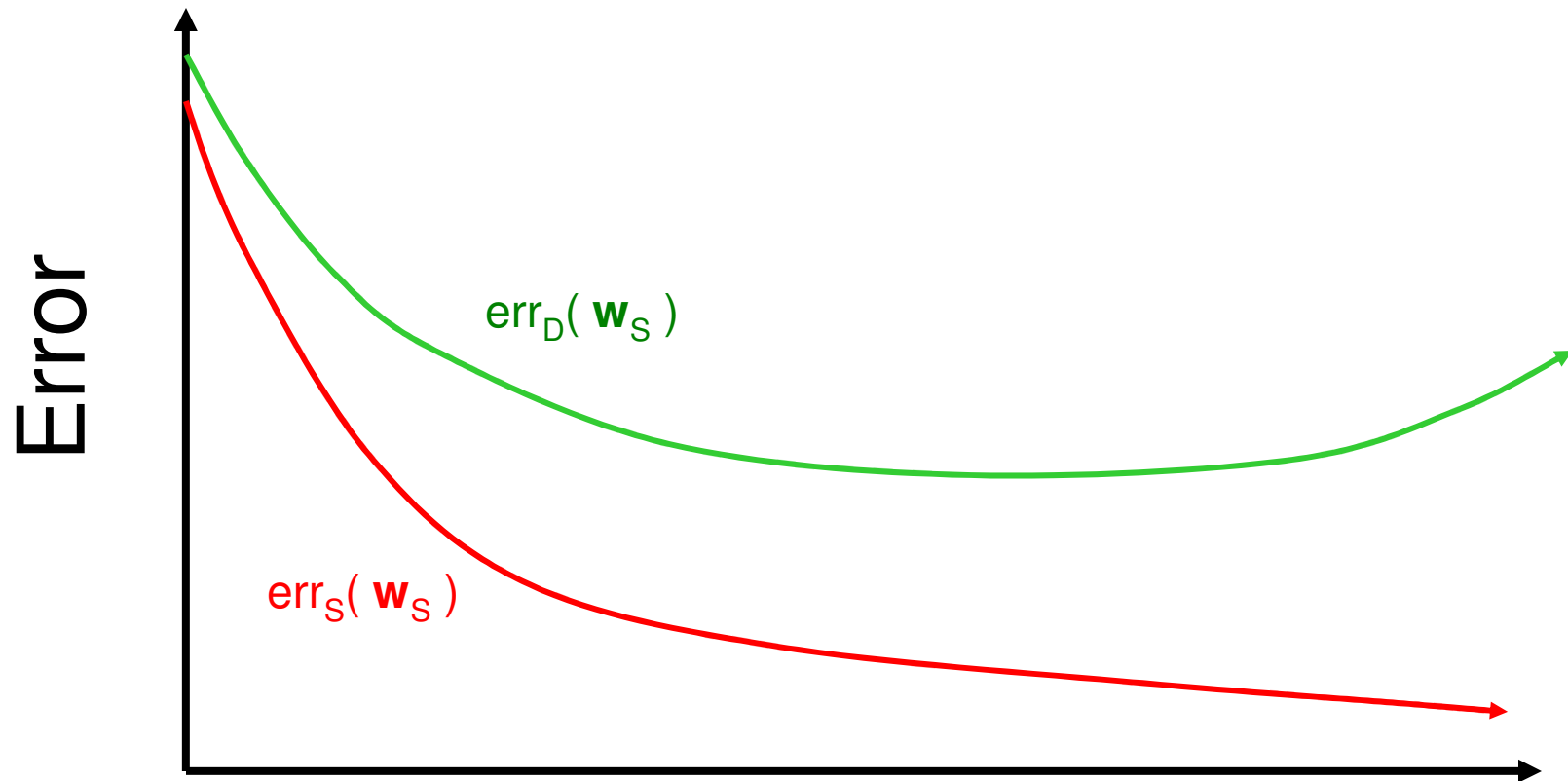
- **WARNING:** Training set error can be poor measure of “quality” of solution
- **Want:** *error over all possible input points, not just training data:*
 - **Prediction error:**

$$\begin{aligned} err_D(\mathbf{w}) &= E_{x,t} \left[\left(t(\mathbf{x}) - \sum_i w_i h_i(\mathbf{x}) \right)^2 \right] \\ &= \int_{x,t} \left(t - \sum_i w_i h_i(\mathbf{x}) \right)^2 D(\mathbf{x}, t) dx dt \end{aligned}$$

Requires $D(\mathbf{x}, t)$ – unknown!

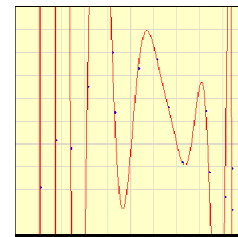


Prediction Error as a function of Model Complexity



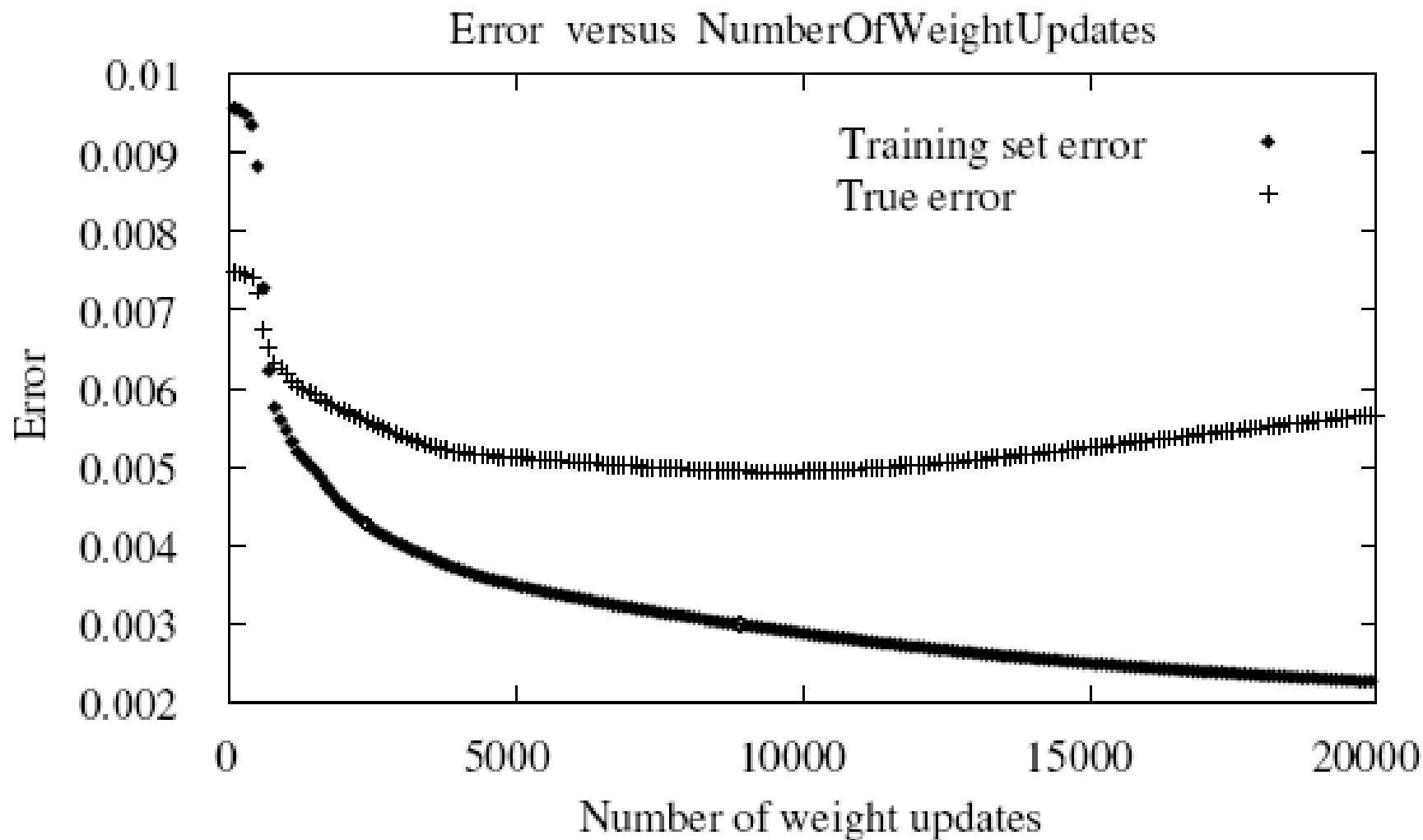
Select points by clicking on the graph or press [Example](#)
Degree of polynomial: FIT to X
 FIT to Y
[Calculate](#) [View Polynomial](#) [Reset](#)

“Model Complexity”
... eg, #basis functions;
degree of poly, ...



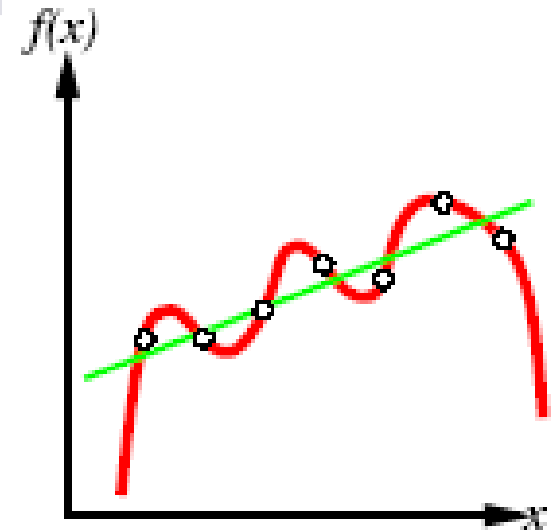
Select points by clicking on the graph or press [Example](#)
Degree of polynomial: FIT to X
 FIT to Y
[Calculate](#) [View Polynomial](#) [Reset](#)

$\underline{\text{err}}_S(\mathbf{w}_S)$ vs $\text{err}_D(\mathbf{w}_S)$



But ...

- $\underline{\text{err}}_S(\mathbf{w}_S) \neq \text{err}_D(\mathbf{w}_S)$
- $\underline{\text{err}}_S(\mathbf{w}_S) \equiv$
Eval \mathbf{w}_S on training set S
 - only approximation to $\text{err}_D(\mathbf{w}_S)$
 - ... can be TOO optimistic!
- “Cheating”
Like being evaluated on test
after seeing SAME test. . .



$$\begin{aligned}\underline{\text{err}}_S(\mathbf{w}) &= 0 \\ \underline{\text{err}}_S(\mathbf{w}) &> 0\end{aligned}$$

Computing Prediction Error

- Computing prediction

$$err_D(\mathbf{w}) = \int_{\mathbf{x}, t} \left(t - \sum_i w_i h_i(\mathbf{x}) \right)^2 D(\mathbf{x}, t) dx dt$$

- Depends on $D(\mathbf{x}, t)$ for every \mathbf{x} – typically not known
- Hard integral

- **New sample**: a set of i.i.d. points

$$S' = \{(\mathbf{x}_1, t_1), \dots, (\mathbf{x}_M, t_M)\} \text{ from } D(\mathbf{x}, t)$$

$$err_D(\mathbf{w}_S) \approx err_{S'}(\mathbf{w}_S) = \frac{1}{|S'|} \sum_{(\mathbf{x}, t) \in S'} \left(t - \sum_i w_i h_i(\mathbf{x}) \right)^2$$

Training Error \neq Prediction Error

- Sampling approximation of prediction error:

$$\text{err}_{S'}(\mathbf{w}_S) \approx \text{err}_D(\mathbf{w}_S)$$

- Training error :

$$\text{err}_S(\mathbf{w}_S) \neq \text{err}_D(\mathbf{w}_S)$$

- Very similar equations!!!

- Why is *training error* a bad measure of *prediction error*?

Training Error \neq Prediction Error

Because you cheated!!!

■ Training error is good estimate for a single \mathbf{w} ,
But you optimized \mathbf{w} with respect to the training error,
and found \mathbf{w} that is good for *this set of instances*

■ **Training error is a (optimistically) biased estimate of prediction error**

■ Very similar equations!!!

□ Why is *training error* a bad measure of *prediction error*?

Test Set Error

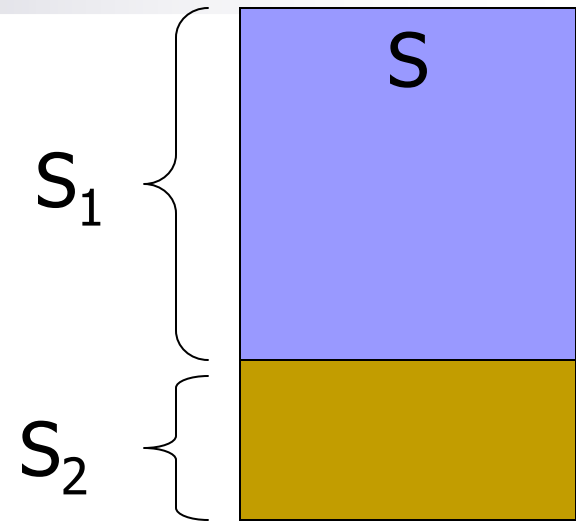
$$\mathbf{w}_S = \mathbf{w}^*(S) = \arg \min_{\mathbf{w}} \sum_{(\mathbf{x}, t) \in S} \left(t - \sum_i w_i h_i(\mathbf{x}) \right)^2$$

- **Randomly** split dataset into two parts:
 - Training data – $S = \{\mathbf{x}_1, \dots, \mathbf{x}_{N_{\text{train}}}\}$
 - Test data – $S' = \{\mathbf{x}_{N_{\text{train}}+1}, \dots, \mathbf{x}_{N_{\text{train}}+N_{\text{test}}}\}$
- Use *training data* to optimize $\mathbf{w} = \mathbf{w}_S$
- **Test set error:**
Given \mathbf{w}_S , estimate error using:

$$\text{err}_{S'}(\mathbf{w}_S) = \frac{1}{|S'|} \sum_{(\mathbf{x}, t) \in S'} \left(t - \sum_i w_i h_i(\mathbf{x}) \right)^2$$

Estimating Error: Hold-Out Set

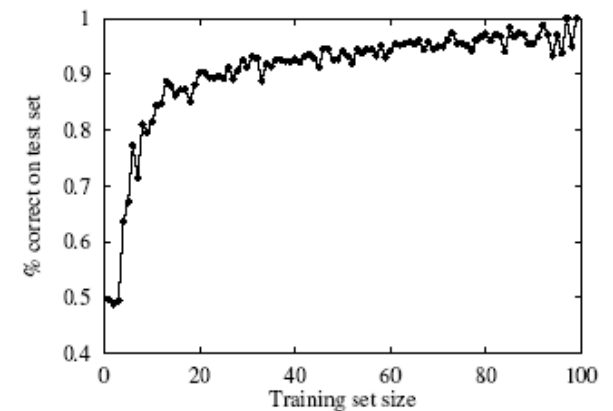
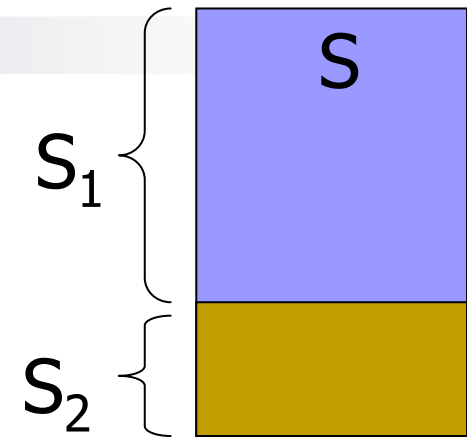
- Run learner $L(\cdot)$ on S
 - produce regressor $w_S = L(S)$What is true error $\text{err}_D(w_S)$?
- Want to return $[w_S, \text{err}_D(w_S)]$
 - ... or at least $[w_S, e]$ where $e \approx \text{err}_D(w_S)$



- Divide S into disjoint S_1, S_2
 - *Train* on S_1 : computing $w_{S_1} := L(S_1)$
 - *Test* on S_2 : $\underline{\text{err}}_{S_2}(w_{S_1})$Return $[w_S, \underline{\text{err}}_{S_2}(w_{S_1})]$
- Why is $\underline{\text{err}}_{S_2}(w_{S_1}) \approx \text{err}_D(w_S)$?
 - As $S_1 \approx S$, $w_{S_1} = L(S_1) \approx L(S) = w_S$
 - $\underline{\text{err}}_{S_2}(w_{S_1})$ is estimate of $\text{err}_D(w_{S_1}) \approx \text{err}_D(w_S)$

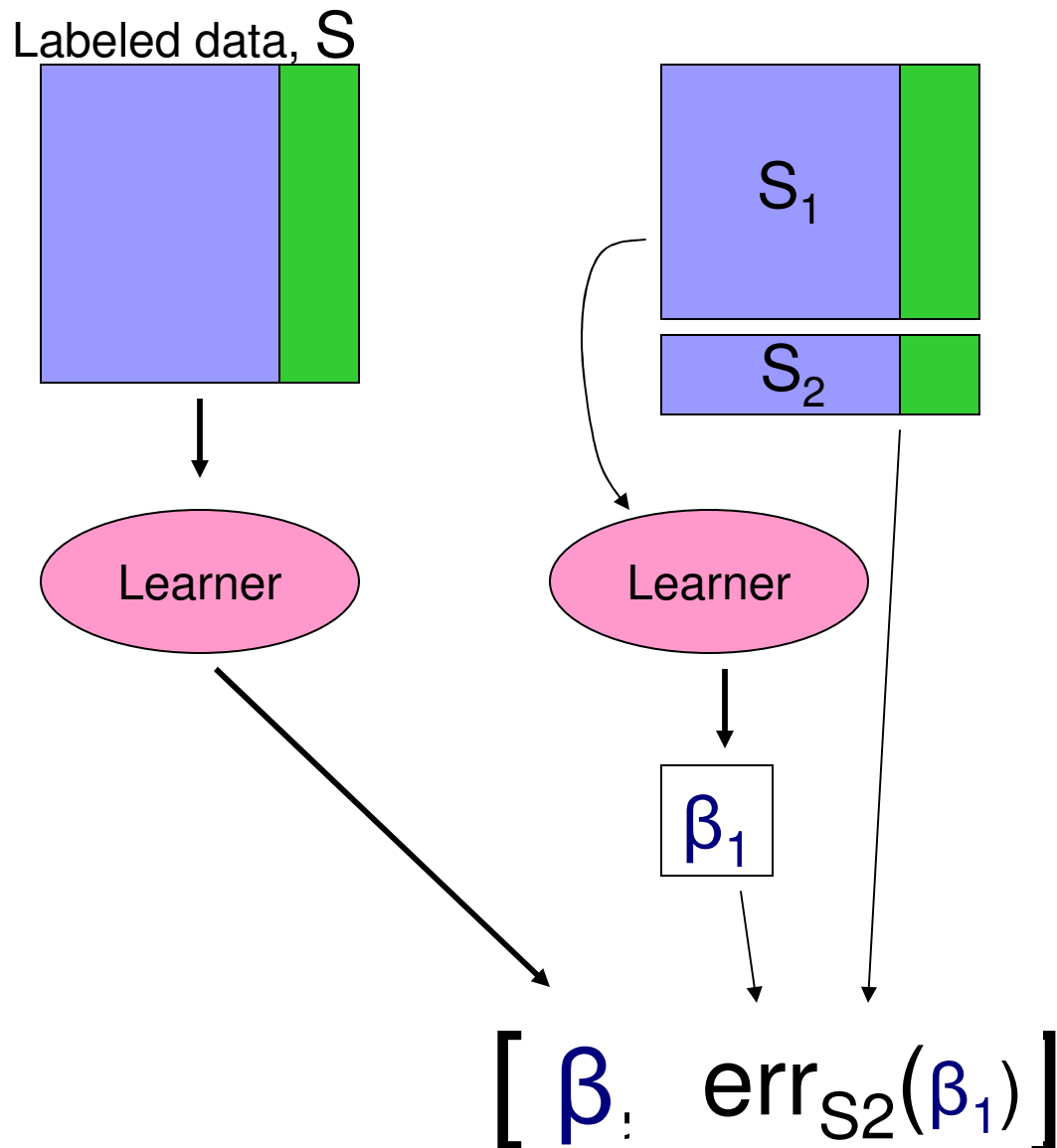
Challenge wrt Hold-Out Set

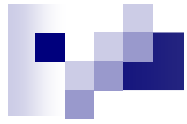
- How to divide S into disjoint S_1, S_2
- As $|S_1| < |S|$,
 $L(S_1)$ not as good as $L(S)$
Learning curve: $L(S)$ improves as $|S|$ increases)
 \Rightarrow want S_1 to be large
- $\underline{err}_{S_2}(w_{S_1})$ is estimate of $err_D(w_S)$
Estimate improves as S_2 gets larger
 \Rightarrow want S_2 to be as large as possible
- As $S = S_1 \cup S_2$, must trade off
quality of classifier $w_{S_1} = L(S_1)$
with
accuracy of estimate $\underline{err}_{S_2}(w_{S_1})$



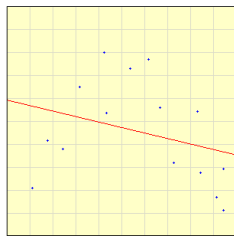
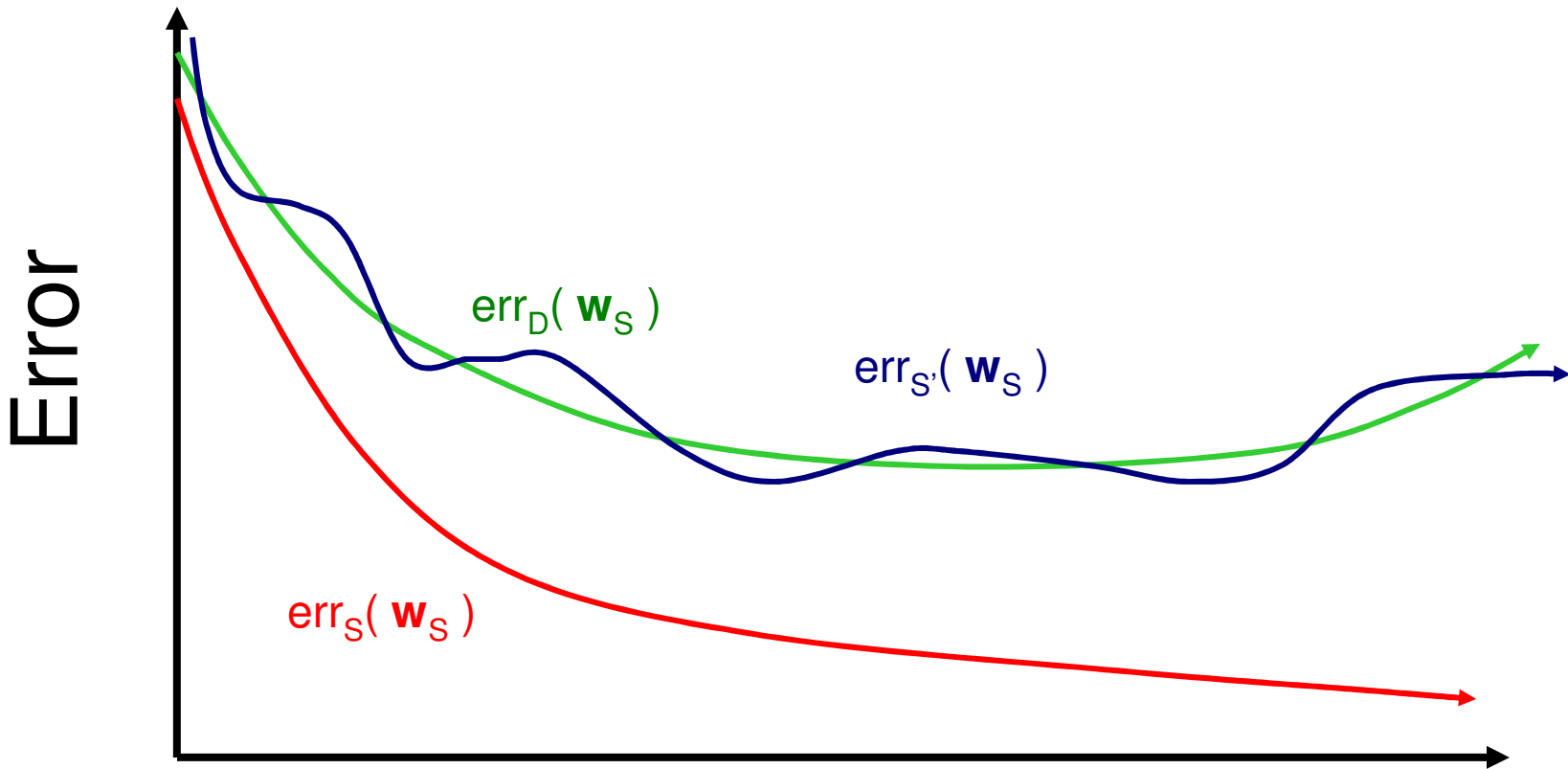
$$|err_S(h) - err_D(h)| \approx \frac{\alpha}{\sqrt{|S|}}$$

Return: [Predictor + Est Quality]



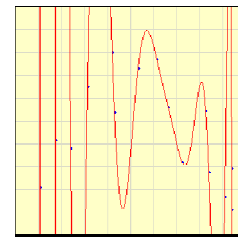


Test Set Error as a function of Model Complexity



Select points by clicking on the graph or press [Example](#)
Degree of polynomial: FIT Y to X
 FIT X to Y
[Calculate](#) [View Polynomial](#) [Reset](#)

“Model Complexity”
... eg, #basis functions;
degree of poly, ...



Select points by clicking on the graph or press [Example](#)
Degree of polynomial: FIT Y to X
 FIT X to Y
[Calculate](#) [View Polynomial](#) [Reset](#)

Estimating Error: Cross Validation

■ “Cross-Validation”

CV(data S , alg L , int k)

Divide S into k disjoint sets $\{ S_1, S_2, \dots, S_k \}$

For $i = 1..k$ do

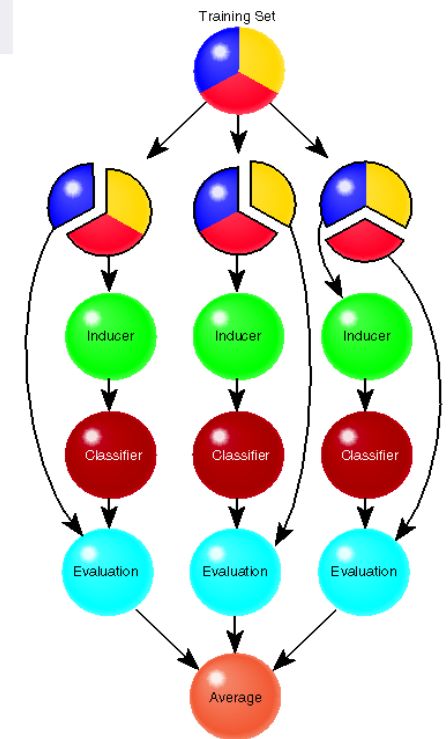
Run L on $S_{-i} = S - S_i$

obtain $h_i := L(S_{-i})$

Evaluate h_i on S_i

$$\text{err}_{S_i}(h_i) = 1/|S_i| \sum_{\langle x, t \rangle \in S_i} [h_i(x) - t]^2$$

Return Average $1/k \sum_i \text{err}_{S_i}(h_i)$



⇒ Less Pessimistic

as train on $(k - 1)/k |S|$ of the data

Comments on Cross-Validation

- Every point used as
 Test 1 time, Training $k - 1$ times
- Computational cost for k -fold Cross-validation ... linear in k
- Should use “balanced CV”
 If class c_i appears in m_i instances,
 insist each S_k include $\approx \frac{1}{k} \frac{m_i}{|S|}$ such instances
- Use **CV(S, L, k)** as ESTIMATE of true error of **L(S)**
 Return **L(S)** and **CV(S, L, k)**
- Leave-One-Out-Cross-Validation $k = m$!
 - eg, for Nearest-Neighbor
- Notice different folds are correlated
 as training sets overlap: $(k-2)/k$ unless $k=2$
- 5×2 -CV
 - Run 2-fold CV, 5 times. . .

Can use CV to estimate parameters in general!



To Form k *Balanced* Folds

1. Partition the data S based on the class:

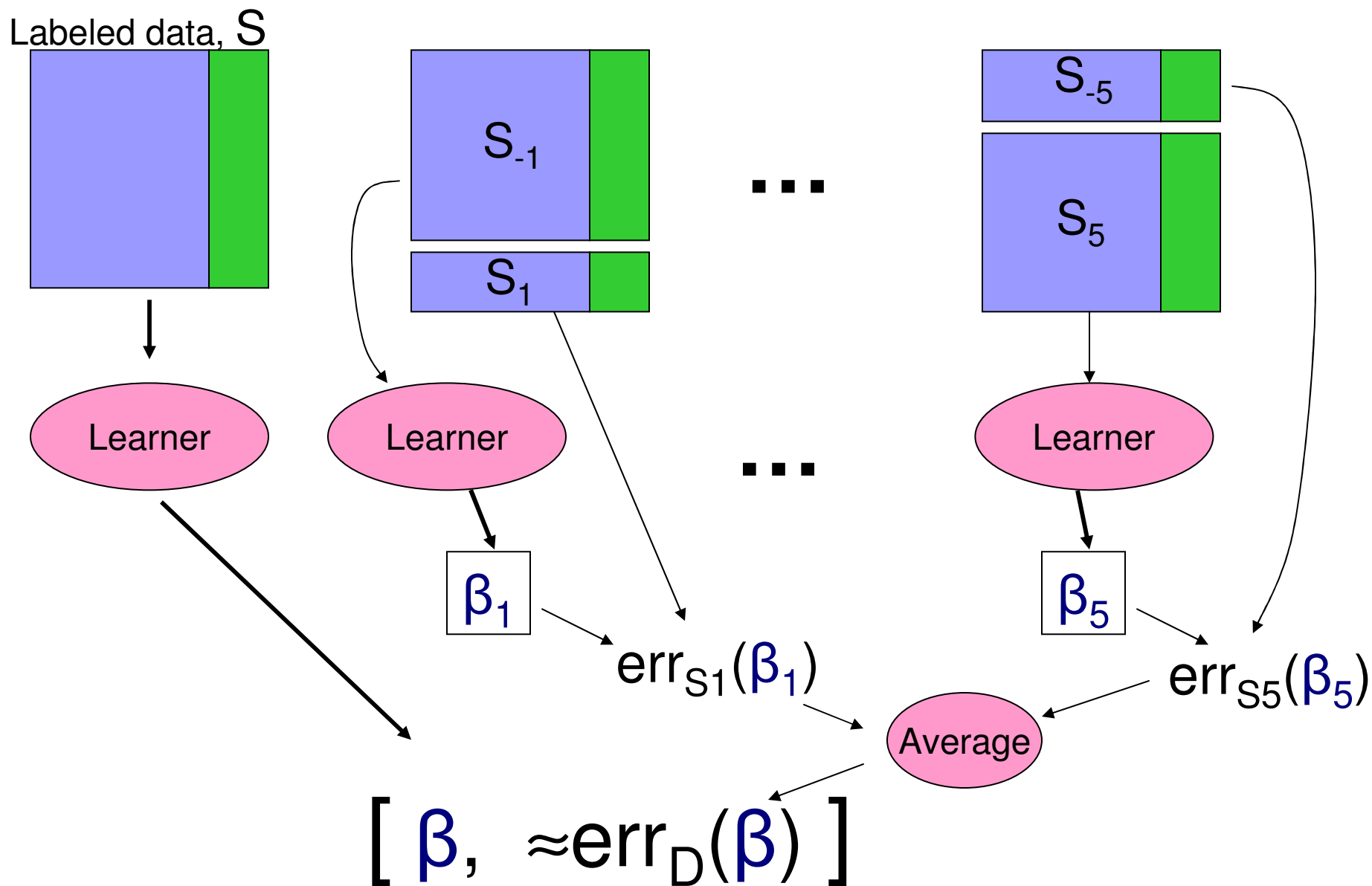
- subset S_+ has all the positive instances,
- subset S_- has all the negative instances.

2. Randomly partition each subset into k folds:

$$S_+ = U \{ S_{+1}, \dots, S_{+k} \}$$

3. $S_j = S_{+j} \cup S_{-j}$ for $j=1..k$

Return: [Predictor + Est Quality]



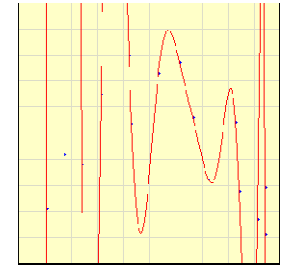
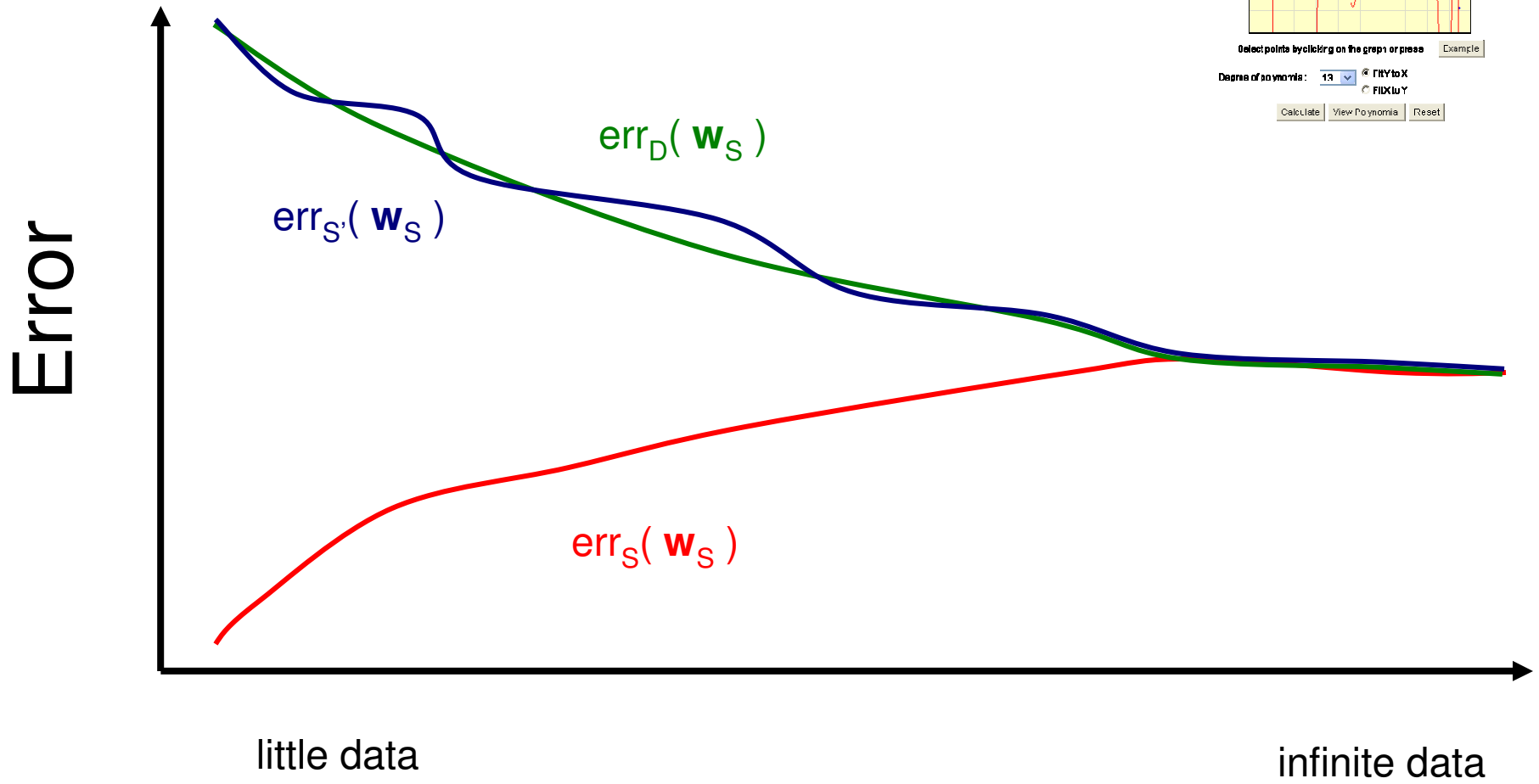
How many points needed for training/testing?

- Very hard question to answer!
 - Too few training points, learned \mathbf{w} is bad
 - Too few test points, you never know if you reached a good solution
- Bounds, such as Hoeffding's inequality can help:

$$P(|\hat{\theta} - \theta^*| \geq \epsilon) \leq 2e^{-2N\epsilon^2}$$

- More on this later this semester, but still hard to answer
- Typically:
 - if you have a reasonable amount of data, pick test set “large enough” for a “reasonable” estimate of error, and use the rest for learning
 - if you have little data, then you need to pull out the big guns...
 - e.g., bootstrapping

Error as function of #Training Examples, for a fixed Model Complexity



Select points by clicking on the graph or press [Example](#)
Degree of polynomial: FIT TO X FIT TO Y

Error Estimators

$$err_D(\mathbf{w}) = \int_{x,t} \left(t - \sum_i w_i h_i(\mathbf{x}) \right)^2 D(\mathbf{x}, t) dx dt$$

Gold Standard!
Unbiased

Uses TRAIN data
... **optimistic**

$$err_S(\mathbf{w}_S) = \frac{1}{|S|} \sum_{(\mathbf{x}, t) \in S} \left(t - \sum_i w_i h_i(\mathbf{x}) \right)^2$$

Approx truth...
Unbiased
... if you are careful

$$err_{S'}(\mathbf{w}_S) = \frac{1}{|S'|} \sum_{(\mathbf{x}, t) \in S'} \left(t - \sum_i w_i h_i(\mathbf{x}) \right)^2$$

Error Estimators

Be careful!!!

Test set only unbiased if you *never never never never* do *any any any* learning/adjustment/... on the test data

Eg,

if you use the test set to select the degree of the polynomial...
no longer unbiased!!!

(We will address this problem later in the semester)

$$\text{err}_{S'}(\mathbf{w}_S) = \frac{1}{|S'|} \sum_{(\mathbf{x}, t) \in S'} \left(t - \sum_i w_i h_i(\mathbf{x}) \right)^2$$

... if you are careful

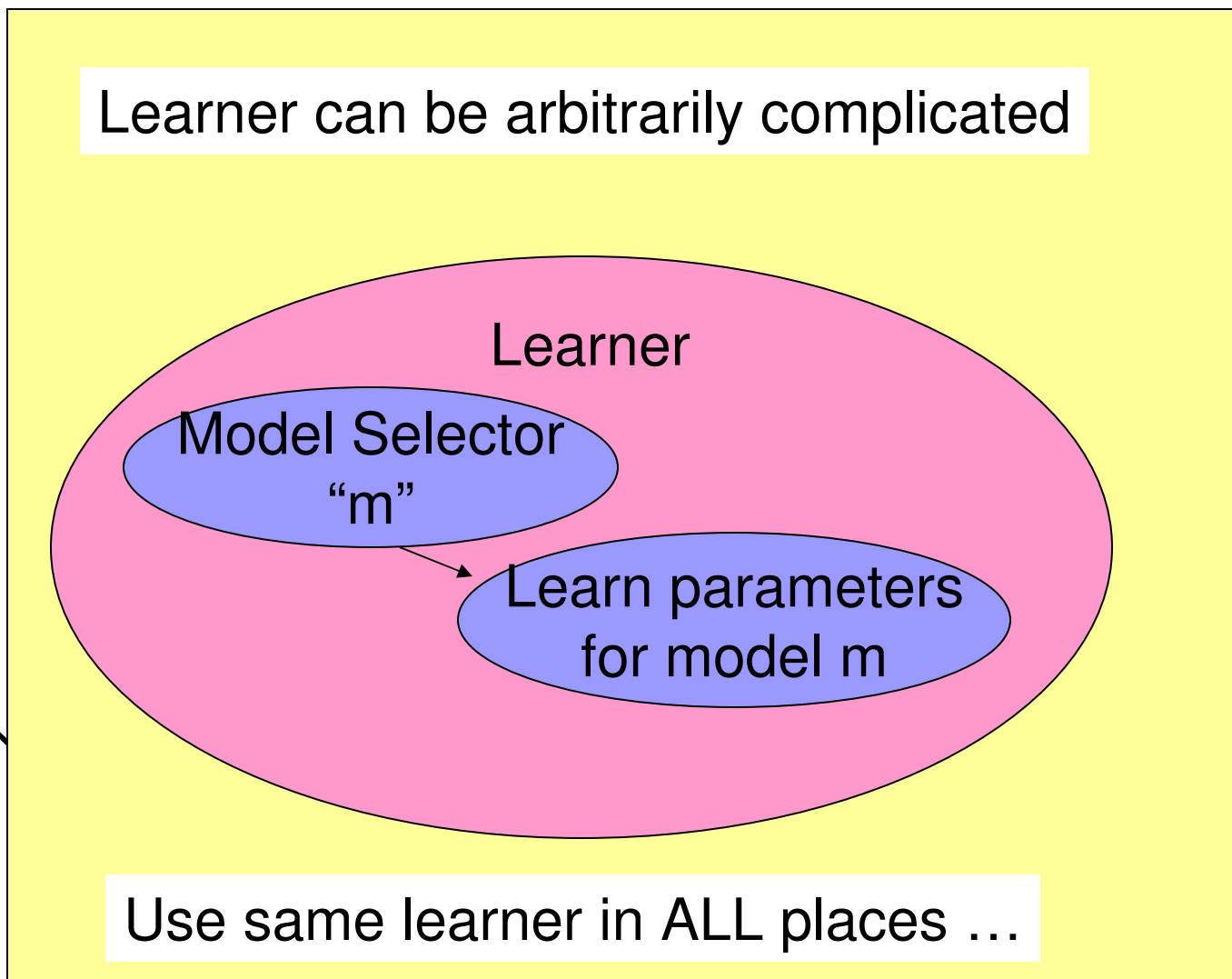
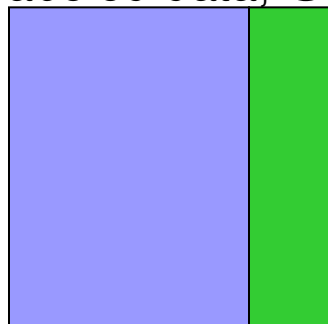


Finding Best Parameters

- Want to learn what “parameters” work best?
 - Best model (RBF vs linear? degree of polynomial)?
Feature selection? Trade-off parameter?...
 - $\operatorname{argmin}_v \{ \operatorname{err}_D(L(S, v)) \}$
- #1?: Try each value on entire dataset.
Report which has smallest TRAINING SET error?
 - $\operatorname{argmin}_v \{ \operatorname{err}_S(L(S, v)) \}$
- #2?: For each value, run 5-fold C-V (wrt entire dataset)
 - $v^* = \operatorname{argmax}_v \{ E[\operatorname{err}_{S_i} (L(S_{-i}, v))] \}$
- Run cross-validation on “best-value” algorithm

Return: [Predictor + Est Quality]

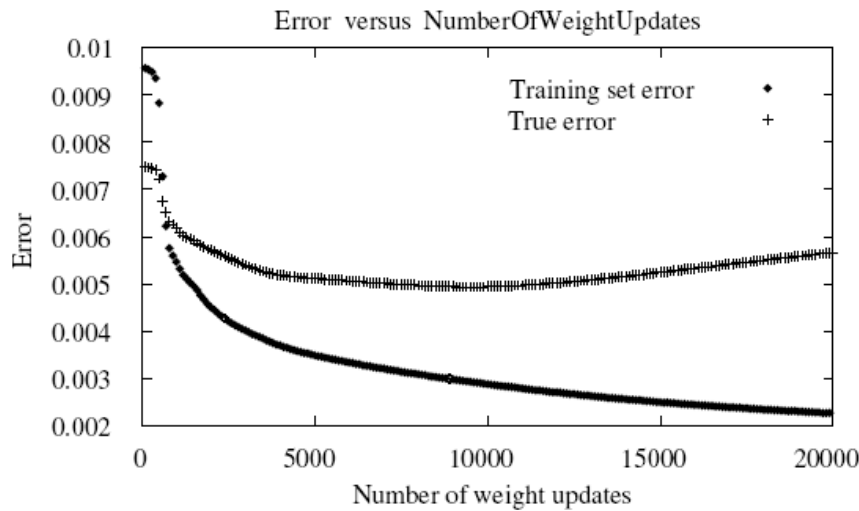
Labeled data, S



$\beta_5)$

[β , $\approx \text{err}_D(\beta)$]

Fit-to-Data \neq Generalization



- h_k = hyp after k updates
 $\widehat{err}_S(h_{20000}) < \widehat{err}_S(h_{10000})$ but
 $err_{\mathcal{D},f}(h_{20000}) > err_{\mathcal{D},f}(h_{10000})$!

■ “Overfitting”

Best “fit-to-data” can find “meaningless regularity” in data (coincidences in the noise)

\Rightarrow bad generalization behavior

- Gen'l: Hypothesis $h \in \mathcal{H}$ overfits training data if \exists alternative hypothesis $h' \in \mathcal{H}$ s.t.

$$\begin{aligned} &\widehat{err}_S(h) < \widehat{err}_S(h') \\ \text{but} \\ &err_{\mathcal{D},f}(h) > err_{\mathcal{D},f}(h') \end{aligned}$$

Summary of Estimating Error

- SetUp: Learner L ,
 - using labeled training data S
 - produces predictor $h_S = L(S)$
- Want $\text{err}_D(h_S)$
 h 's Generalization Error over distribution D
 - to evaluate predictor h_S
 - to decide among possible predictors
 - to evaluate learner
- But depends on $D(\mathbf{x}, t)$:
not known!

Estimating $\text{err}_D(h_S)$

1. Training Set Error
 - Use h_S 's empirical error on S
 $\text{err}_S(h_S)$
⇒ Very Optimistic
2. Hold Out Error
 - Divide $S = S_1 \cup S_2$;
Return $\text{err}_{S_2}(h_{S_1})$
⇒ Slightly Pessimistic
3. Cross Validation
 - $1/k \sum_i \text{err}_{S_i}(L(S_{-i}))$
⇒ Slightly Less Pessimistic

For evaluating GENERAL PREDICTORS

- classifiers, regressors
- ... best values for parameters