

Assessing Performance

CSE 446: Machine Learning
Slides created by Emily Fox (mostly)

April 8, 2019

$$(\vec{x}_i, y_i) \rightarrow (\vec{x}_i - \vec{m}, y_i)$$

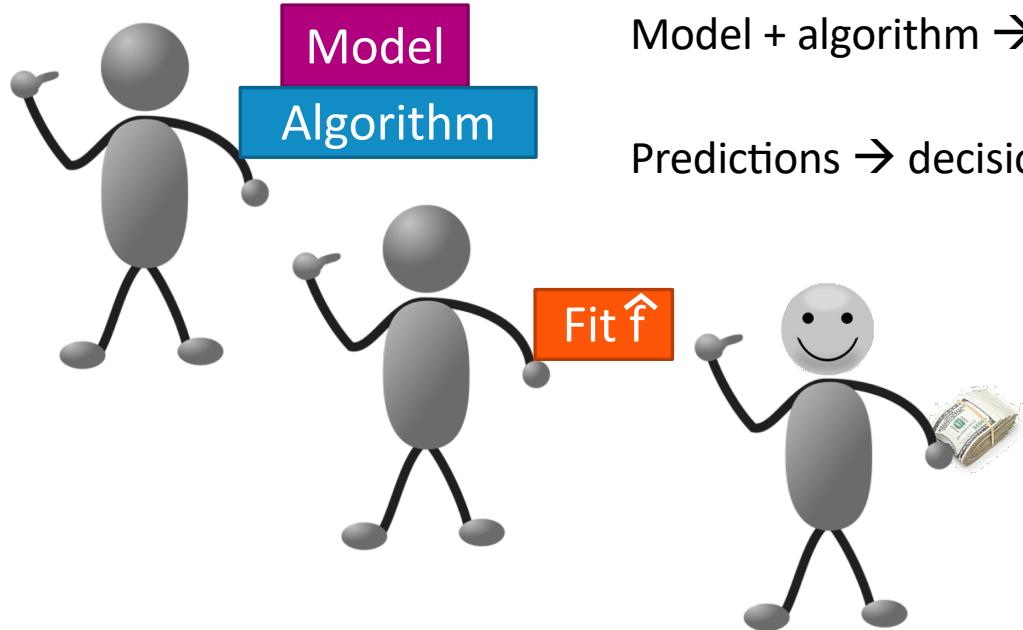
$$\vec{m} = \frac{1}{n} \sum_{i=1}^n \vec{x}_i$$

$$\vec{x}_i = \begin{pmatrix} 1, x_i[1], x_i[2], \dots, x_i[d] \end{pmatrix}$$

To fit these more general functions

- Start with input features $\mathbf{x} = (x[1], x[2], \dots, x[d])$ and training set: $\{(x_i, y_i)\}_{i=1..n}$
- Define feature map that transforms each x_i to higher dimensional feature vector $\mathbf{h}(x_i)$.
- Model: $y_i = \sum_{j=1}^p w_j h_j(x_i) + \varepsilon_i$
- Find $\hat{\mathbf{w}}$ that minimizes $\text{RSS} = \sum_{i=1}^n (y_i - \sum_{j=1}^p w_j h_j(x_i))^2$
 $= (\mathbf{y} - \mathbf{H}\mathbf{w})^\top (\mathbf{y} - \mathbf{H}\mathbf{w})$
- Solution: $\mathbf{H}^\top \mathbf{H} \hat{\mathbf{w}} = \mathbf{H}^\top \mathbf{y}$

How good are my predictions?



Model + algorithm → fitted function

Predictions → decisions → outcome

Measuring loss

(\vec{x}, \vec{y})
output.

Loss function:

$$L(y, f_{\hat{w}}(\vec{x}))$$

actual value
 $\hat{f}(\vec{x}) = \text{predicted value } \hat{y}$

Cost of using \hat{w} at \vec{x}
when y is true

Examples:

Squared error: $L(y, f_{\hat{w}}(\vec{x})) = (y - f_{\hat{w}}(\vec{x}))^2$

Absolute error: $L(y, f_{\hat{w}}(\vec{x})) = |y - f_{\hat{w}}(\vec{x})|$

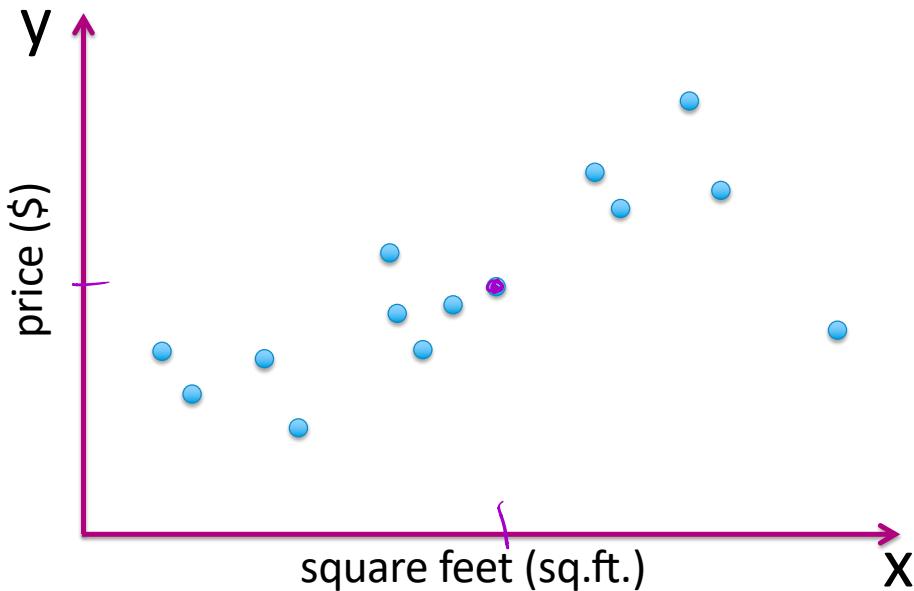
“Remember that all models are wrong; the practical question is how wrong do they have to be to not be useful.” George Box, 1987.

Assessing the loss

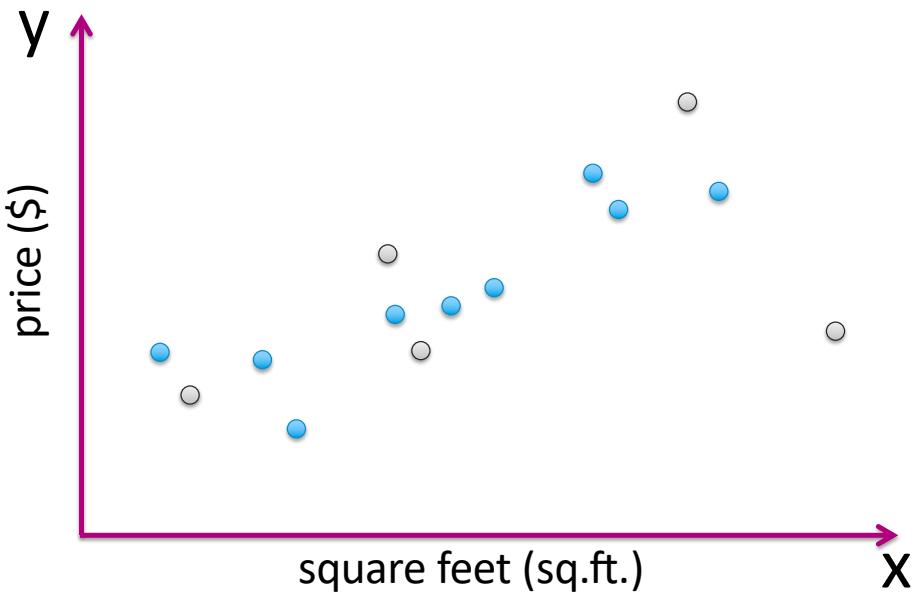
Assessing the loss

Part 1: Training error

Start with a data set

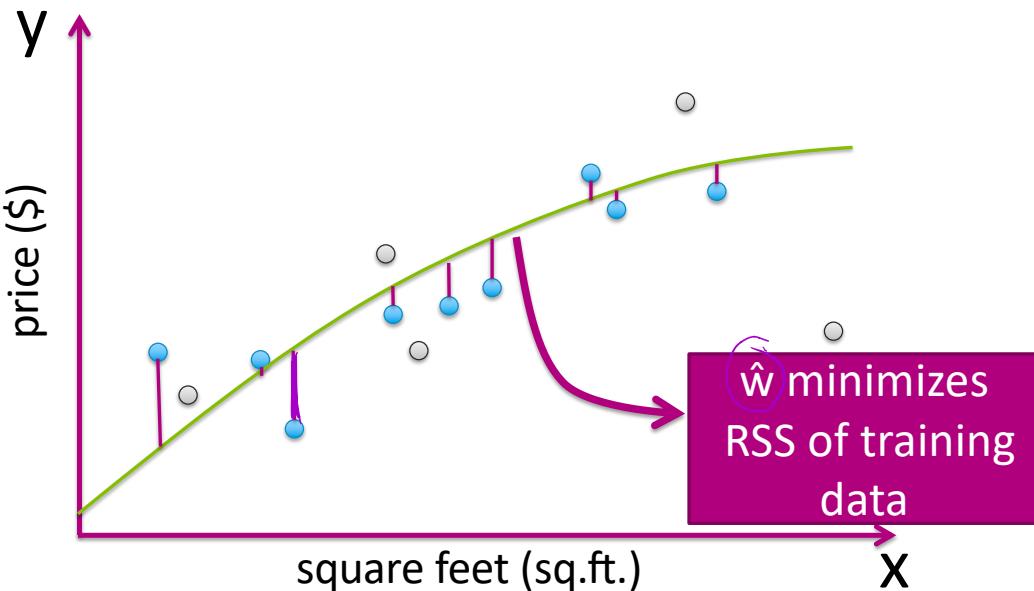


Define training data



$$\{\vec{x}_i, y_i\}_{i=1}^n$$

Example:
Fit quadratic to minimize RSS



Compute training error

1. Define a loss function $L(y, f_{\hat{w}}(x))$
 - E.g., squared error,...
2. Training error using $f_{\hat{w}}$
= avg. loss on houses in training set

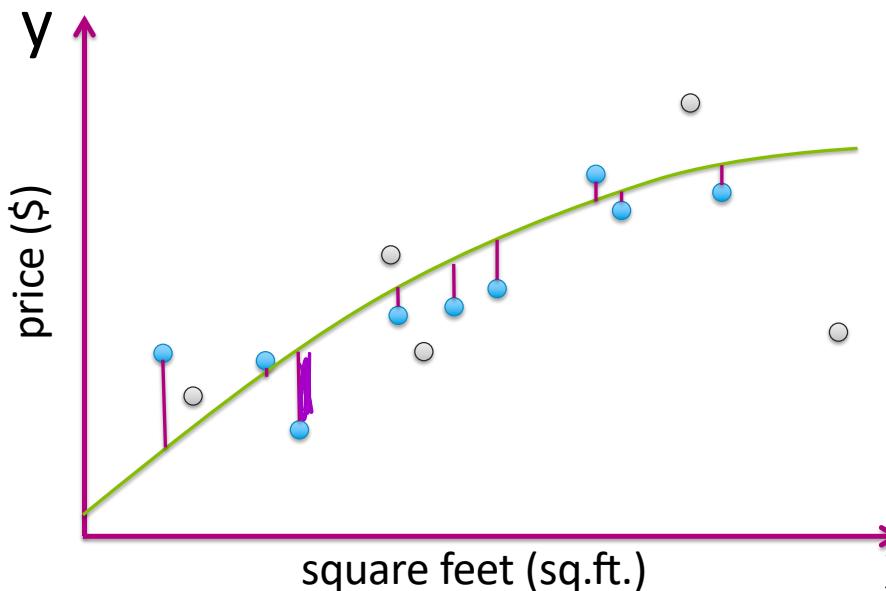
$$= \frac{1}{n} \sum_{i=1}^n L(y_i, f_{\hat{w}}(x_i))$$



fit using training data

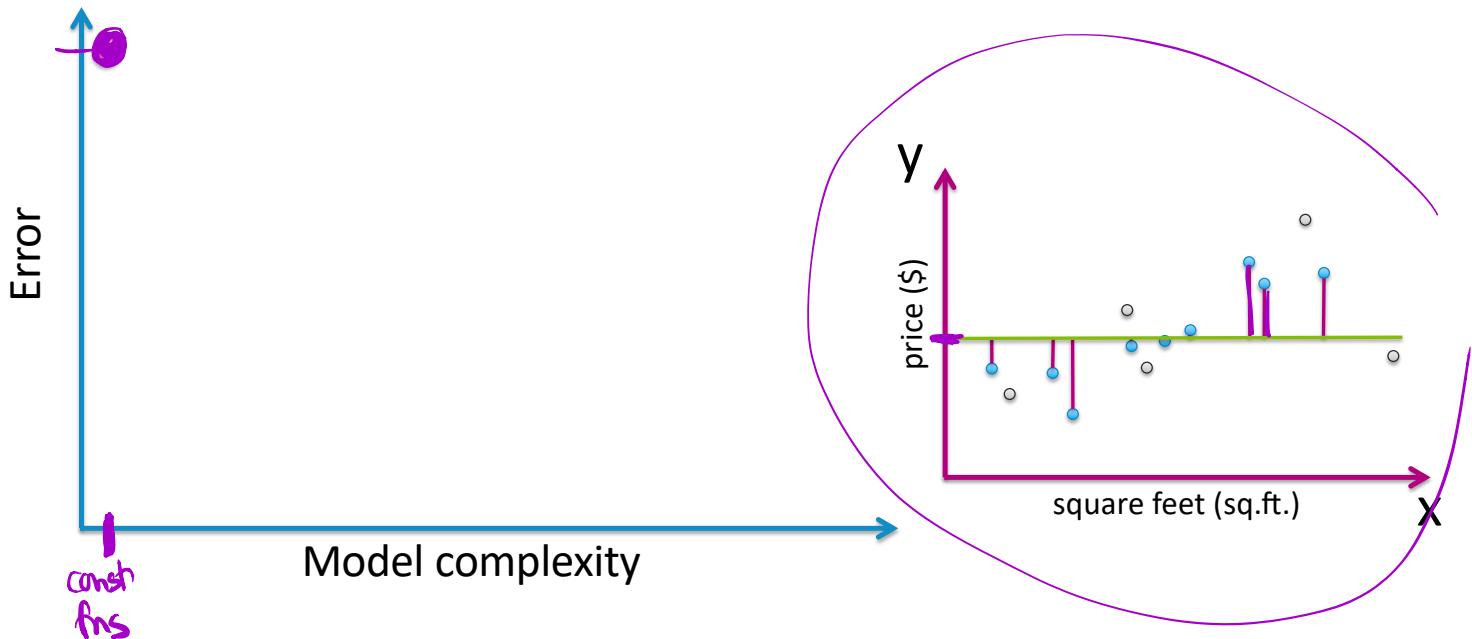
Example:

Use squared error loss $(y - f_{\hat{w}}(x))^2$



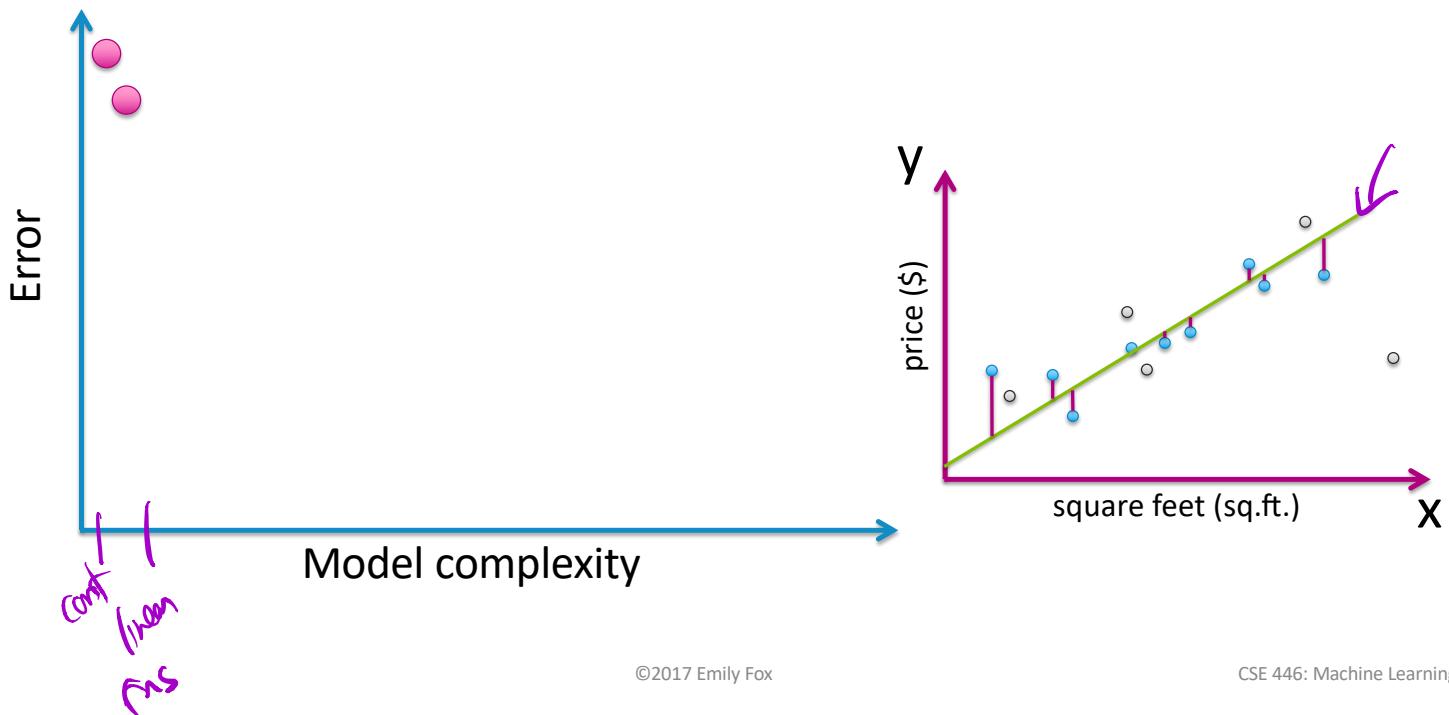
Training error (\hat{w}) = $1/n * [(\$_{train\ 1} - f_{\hat{w}}(sq.ft._{train\ 1}))^2 + (\$_{train\ 2} - f_{\hat{w}}(sq.ft._{train\ 2}))^2 + (\$_{train\ 3} - f_{\hat{w}}(sq.ft._{train\ 3}))^2 + \dots \text{ include all training houses}]$

Training error vs. model complexity

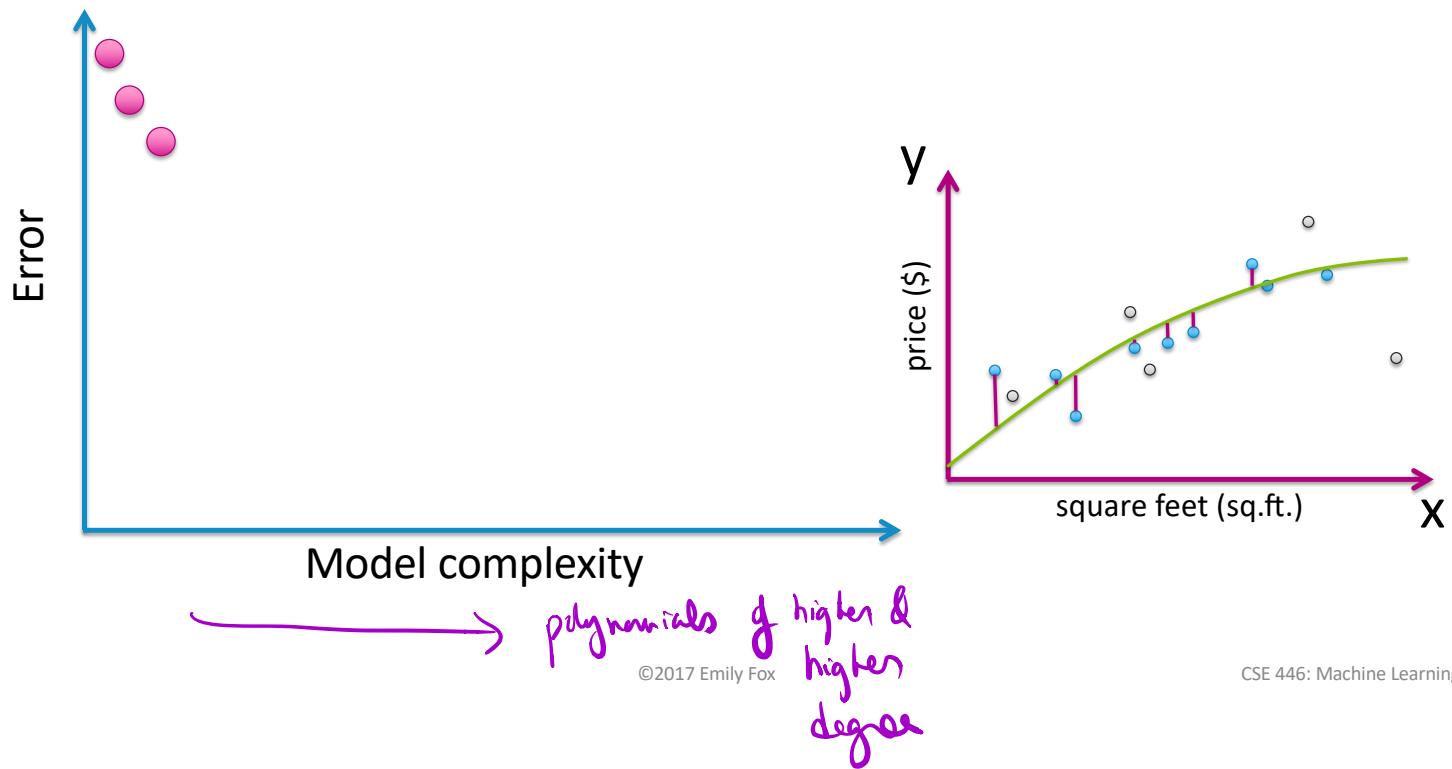


Training error vs. model complexity

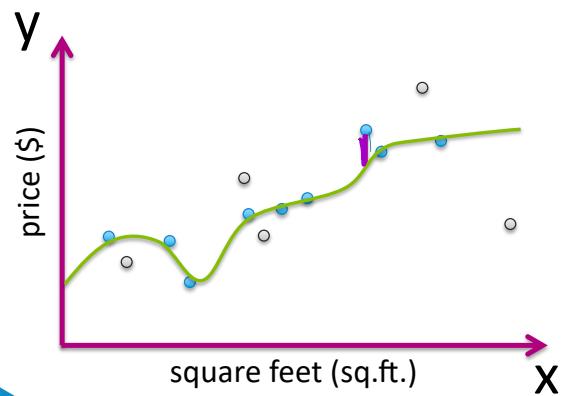
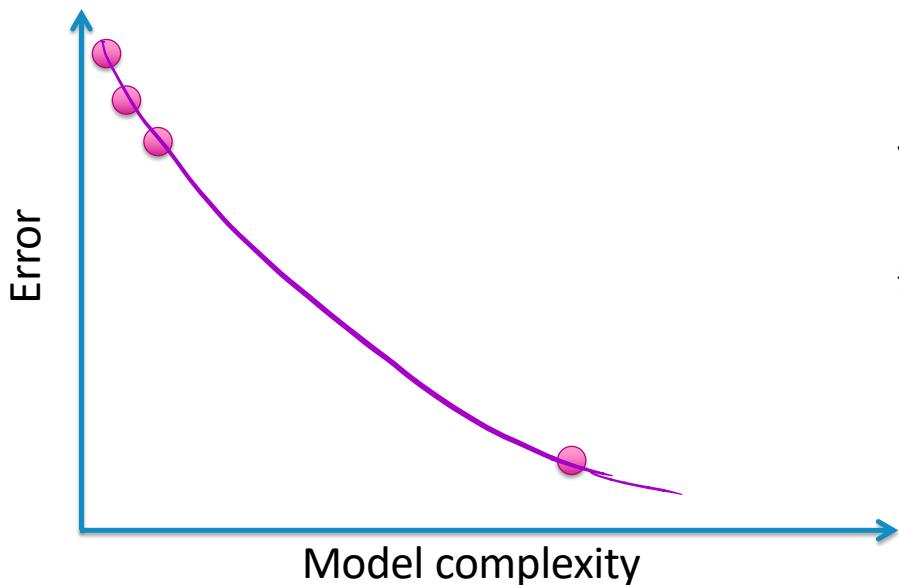
blue pts ~
training set.



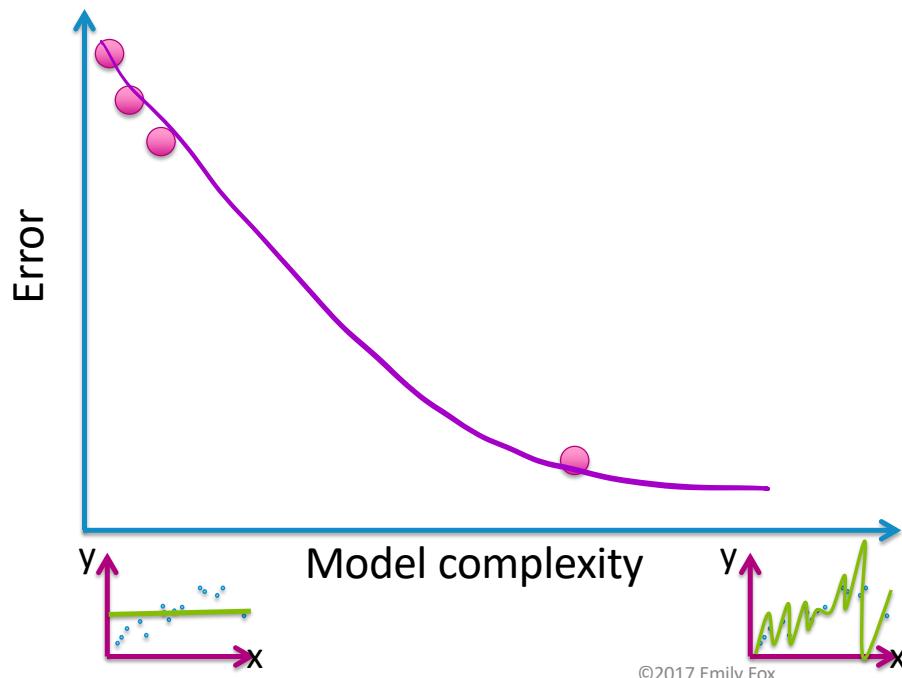
Training error vs. model complexity



Training error vs. model complexity

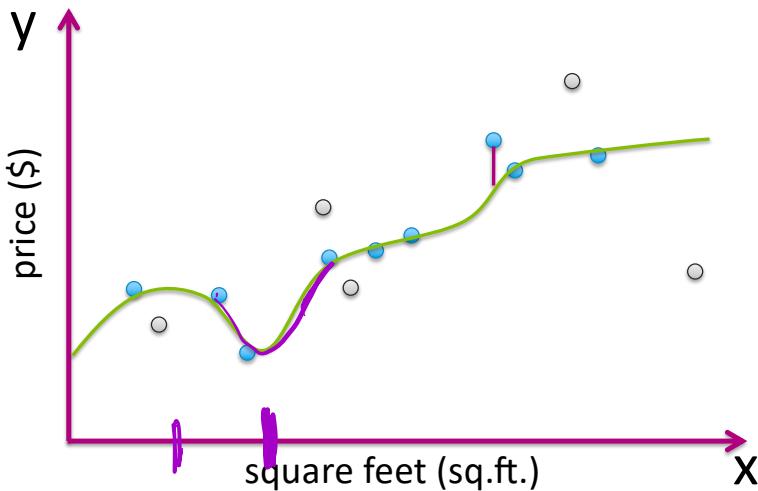


Training error vs. model complexity



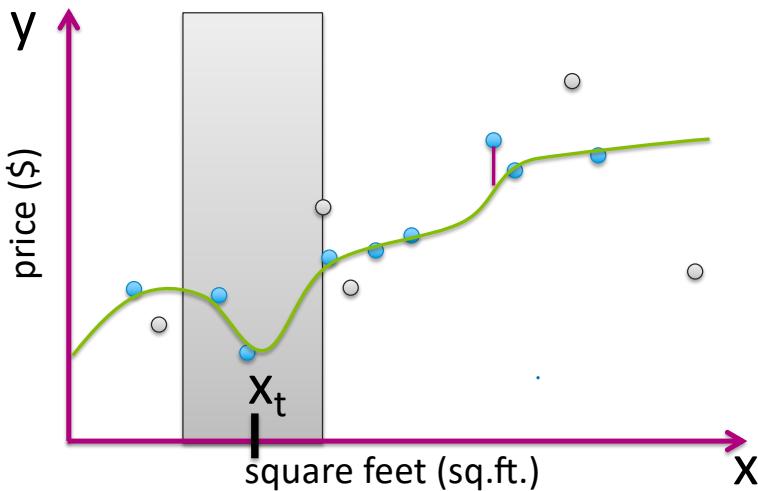
Is training error a good measure of predictive performance?

How do we expect to perform on a new house?



Is training error a good measure of predictive performance?

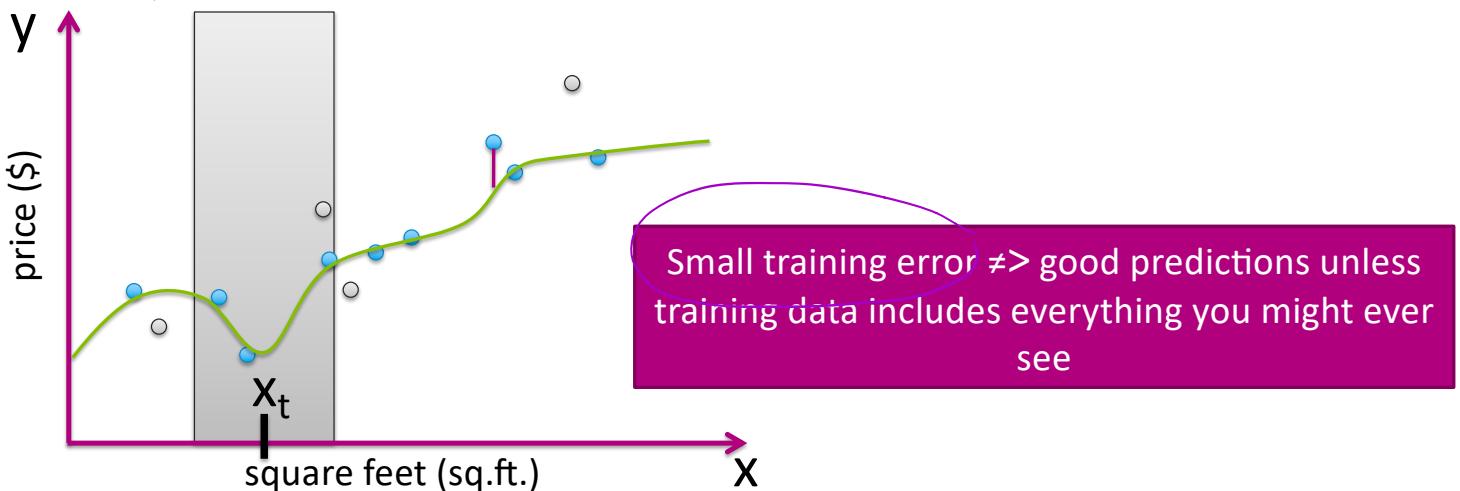
Is there something particularly bad about having x_t sq.ft.??



Is training error a good measure of predictive performance?

Issue:

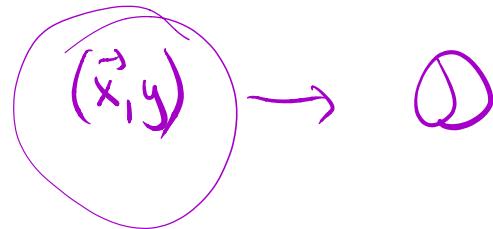
Training error is overly optimistic... \hat{w} was fit to training data



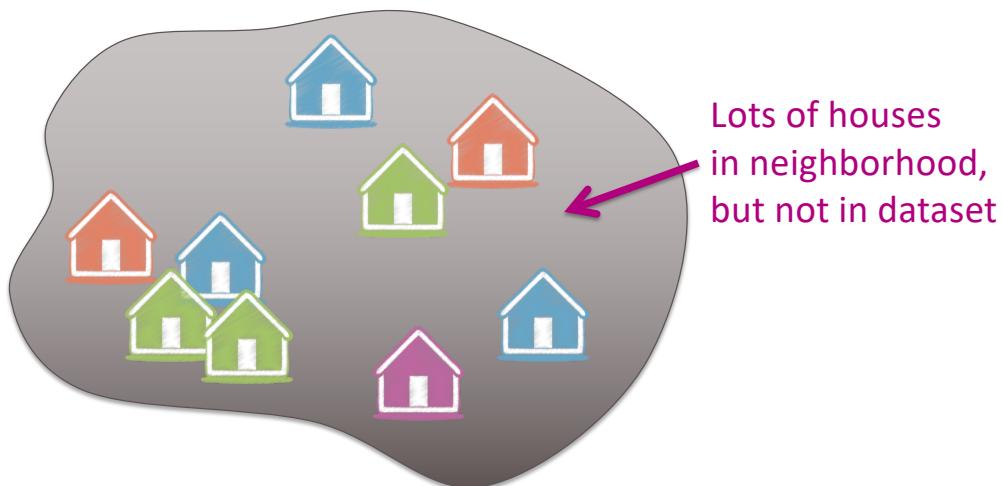
Assessing the loss

Part 2: Generalization (true) error

Generalization error



Really want estimate of loss over all possible (,) pairs

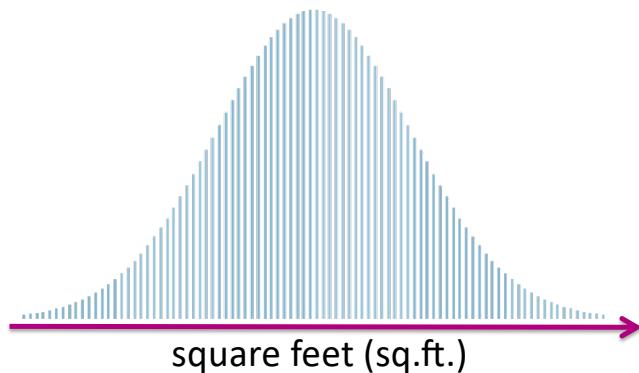


$$p(\vec{x}_i | y)$$

Distribution over houses

$$p(\vec{x})$$

In our neighborhood, houses of what # sq.ft. ()
are we likely to see?



$$(\vec{x}_i | y)$$

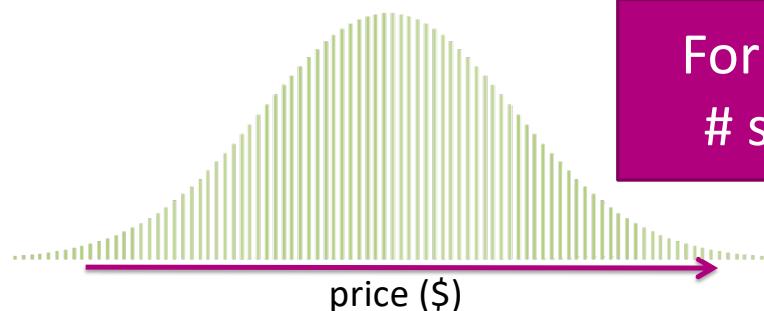
$$p(x,y) = p(x) p(y|x)$$

Distribution over sales prices

For houses with a given # sq.ft. (🏠), what house prices \$ are we likely to see?

$$p(y|x)$$

For fixed
sq.ft.



Generalization error definition

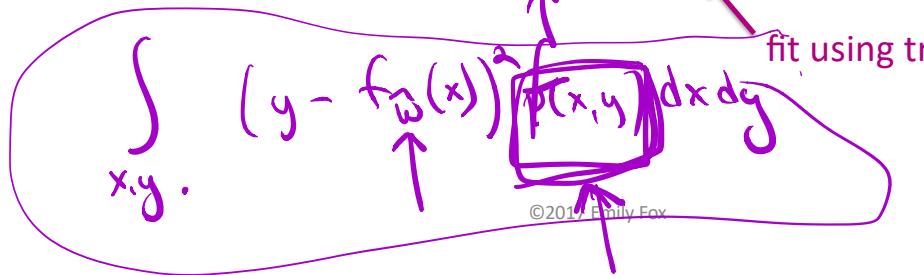
Really want estimate of loss over all possible (,\$) pairs

Formally:

average over all possible
(x,y) pairs weighted by
how likely each is

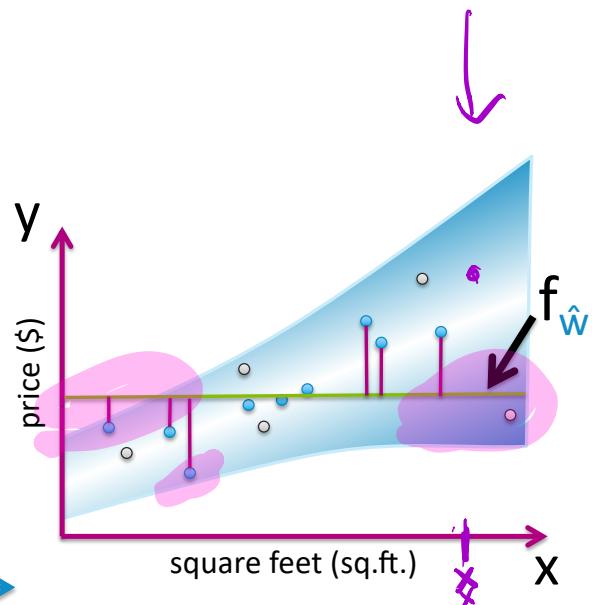
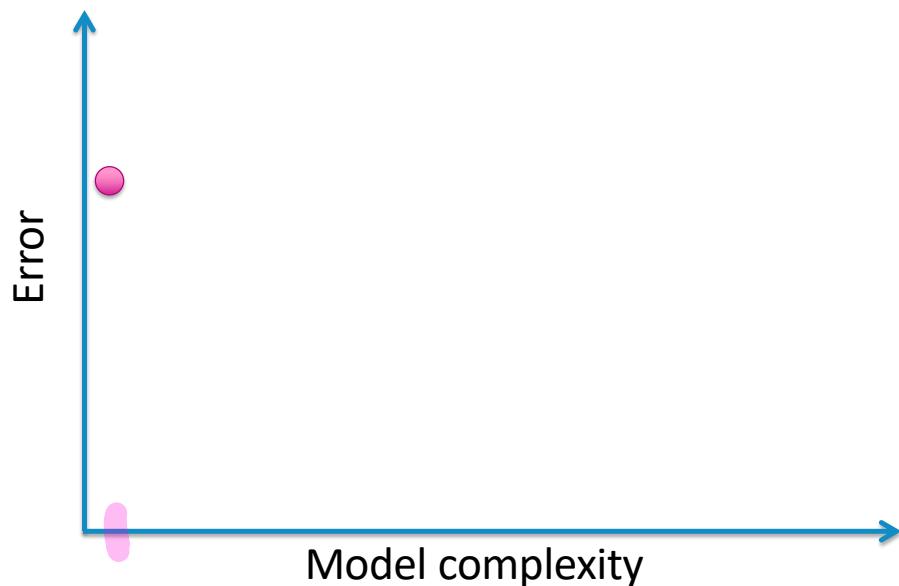


$$\text{generalization error} = E_{x,y} [L(y, f_{\hat{w}}(x))]$$

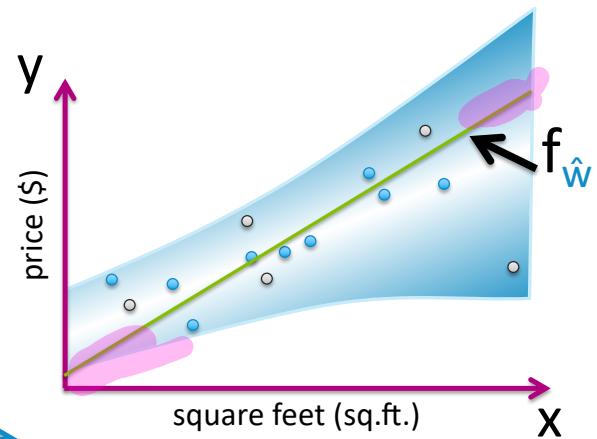
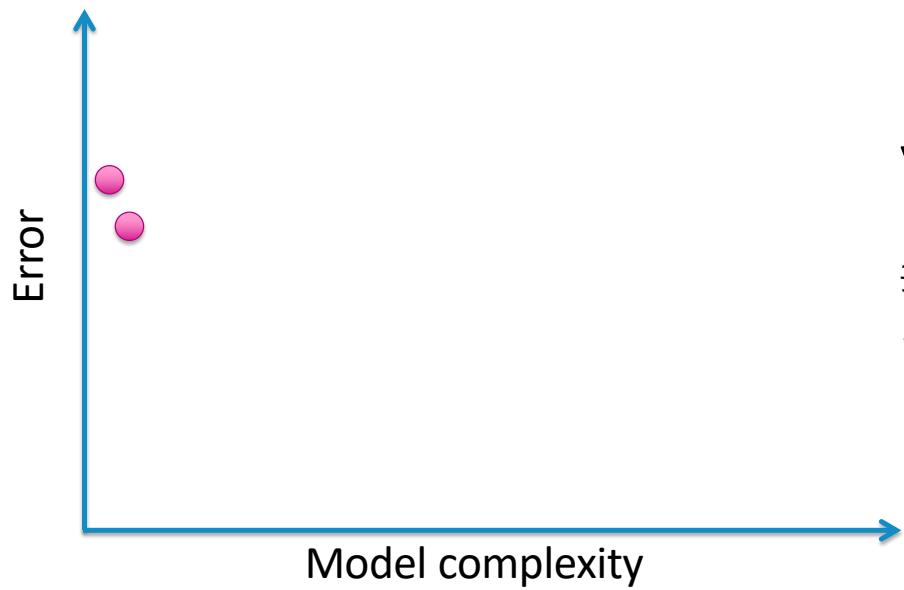


$$p(y|x)$$

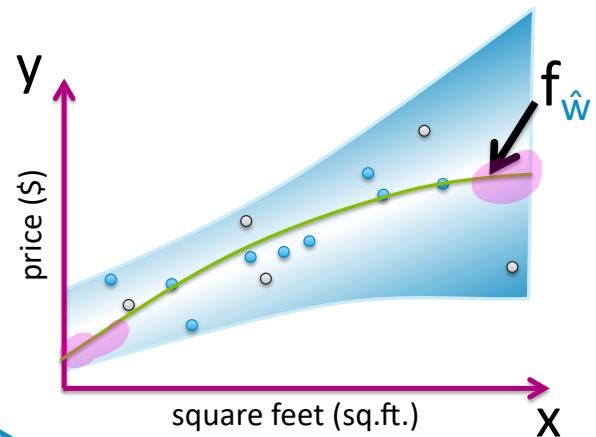
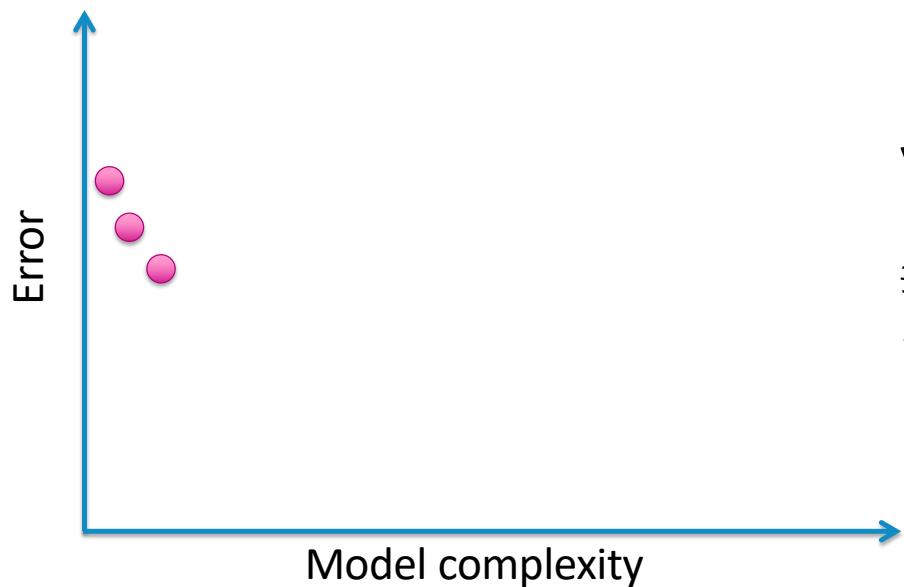
Generalization error vs. model complexity



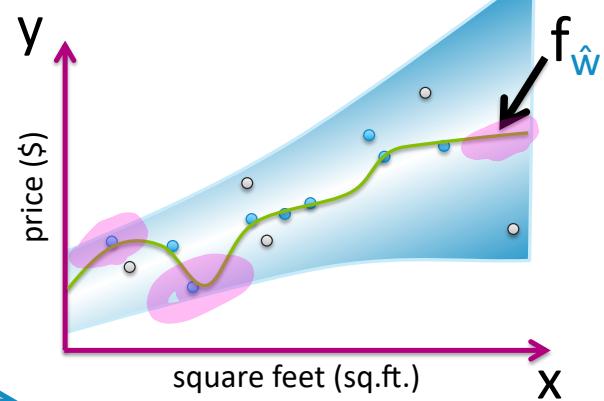
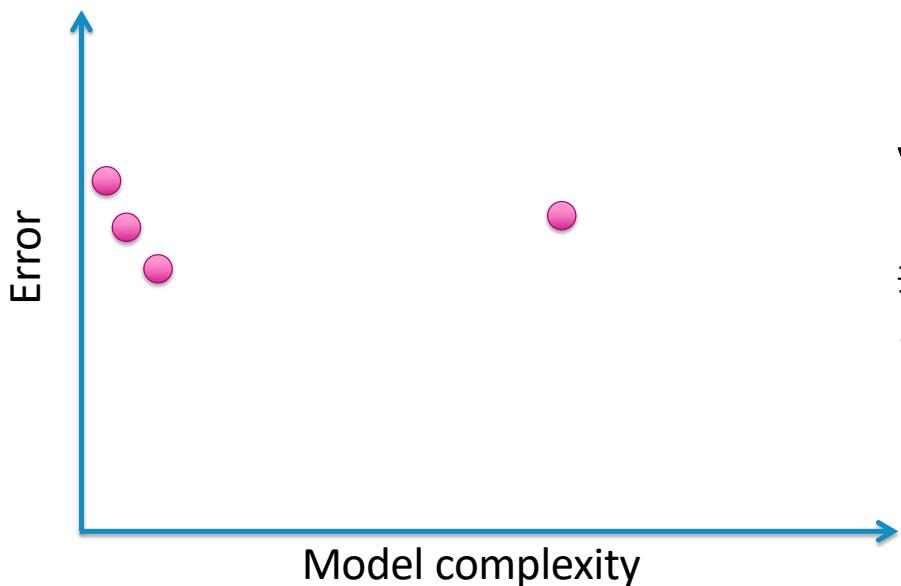
Generalization error vs. model complexity



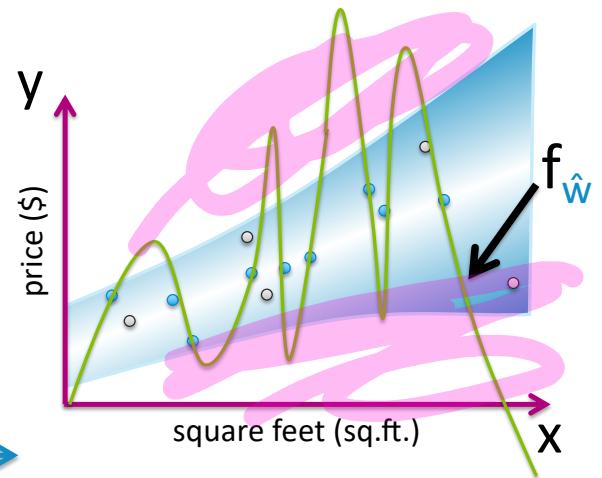
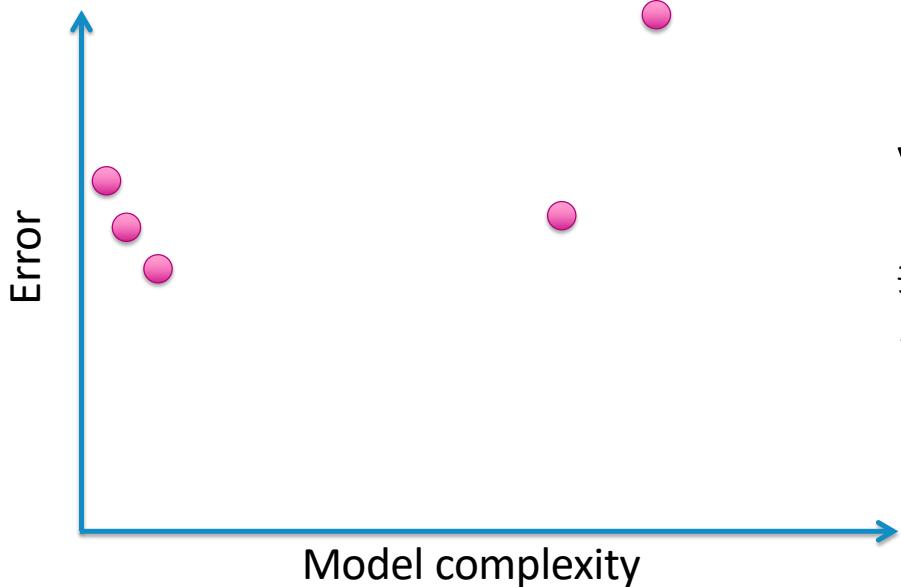
Generalization error vs. model complexity



Generalization error vs. model complexity

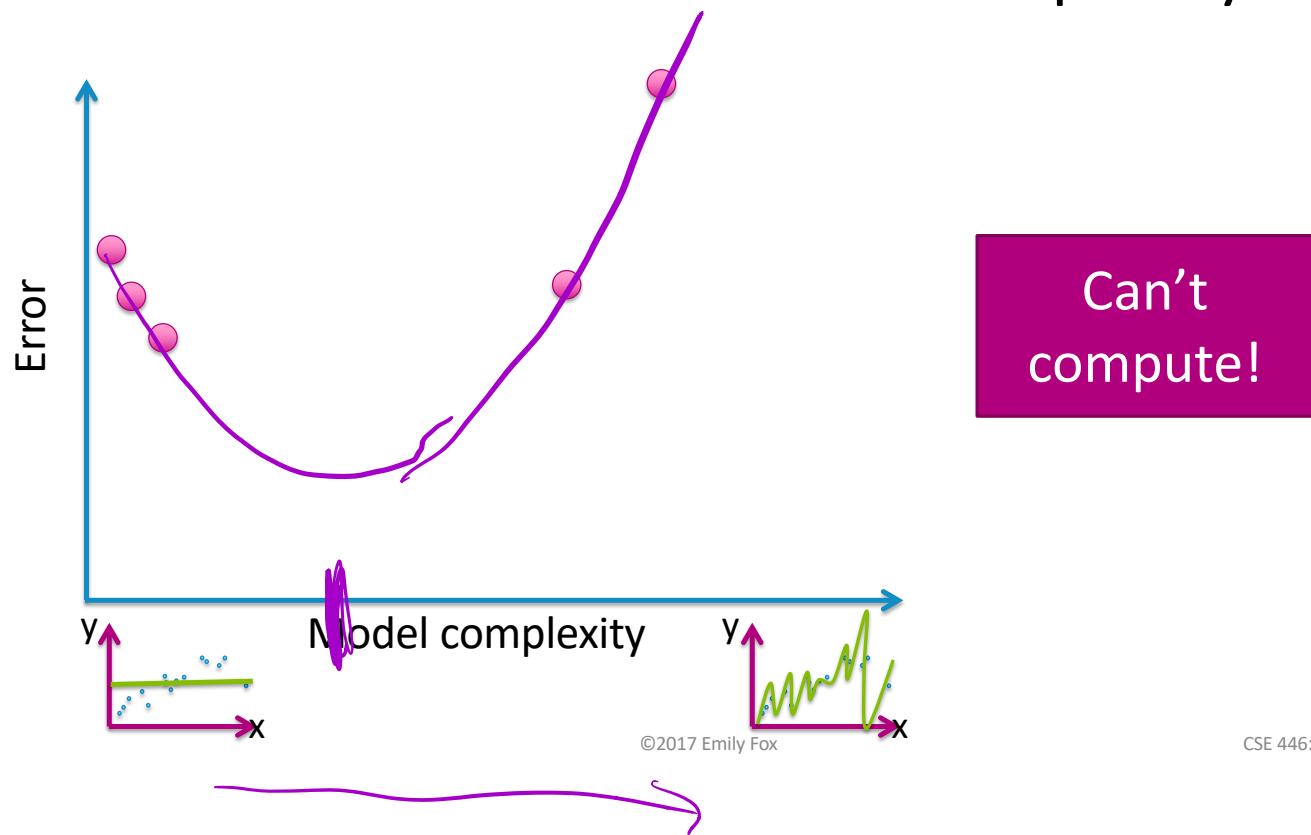


Generalization error vs. model complexity



$$P(x_1, y)$$

Generalization error vs. model complexity



Assessing the loss

Part 3: Test error

Approximating generalization error

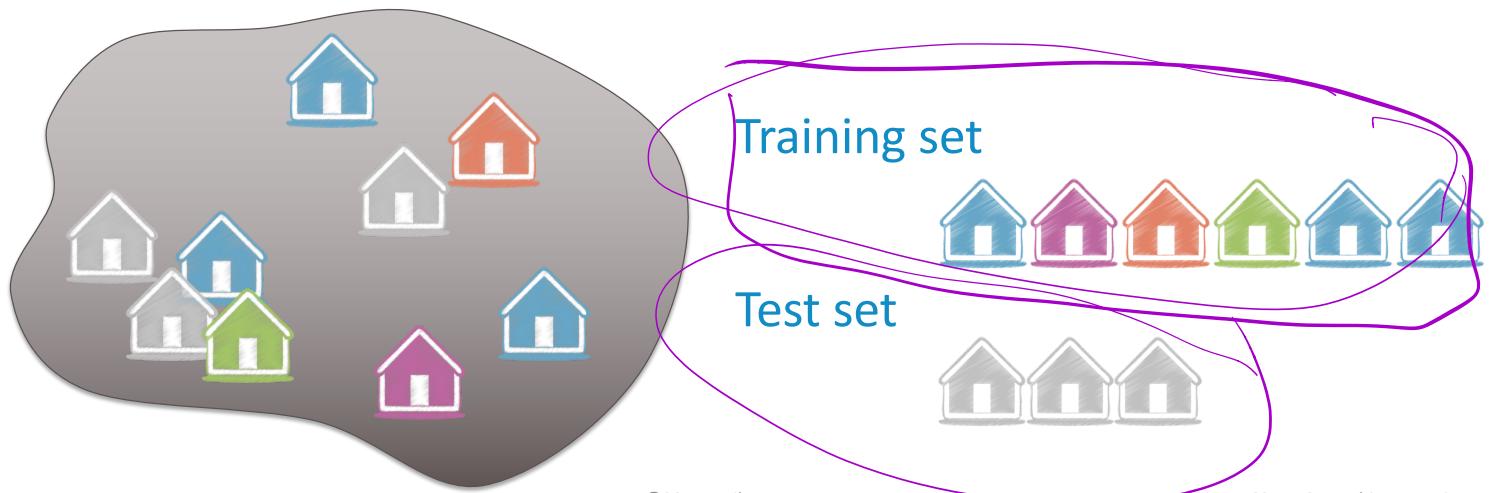
Wanted estimate of loss over all possible (, ) pairs



Approximate by looking at
houses not in training set

Forming a test set

Hold out some (, ) that are *not* used for fitting the model



Forming a test set

Hold out some (, ) that are *not* used for fitting the model



Proxy for “everything you might see”

Test set



every sample i.i.d.
①

Compute test error

Test error

= avg. loss on houses in **test set**

$$= \frac{1}{N_{test}} \sum_{i \text{ in test set}} L(y_i, f_{\hat{w}}(x_i))$$

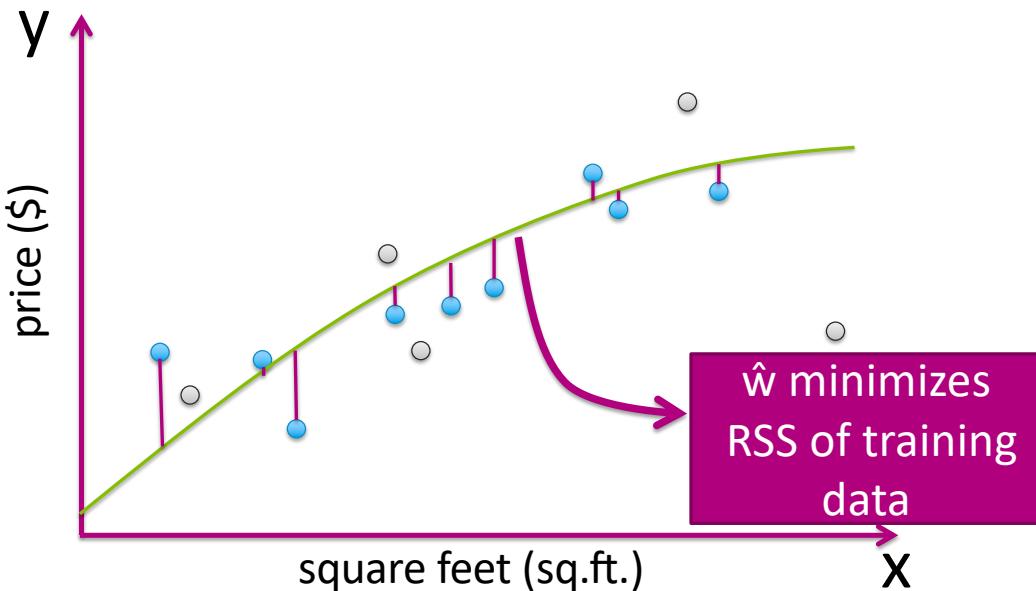
test points

fit using **training data**

has never seen
test data!

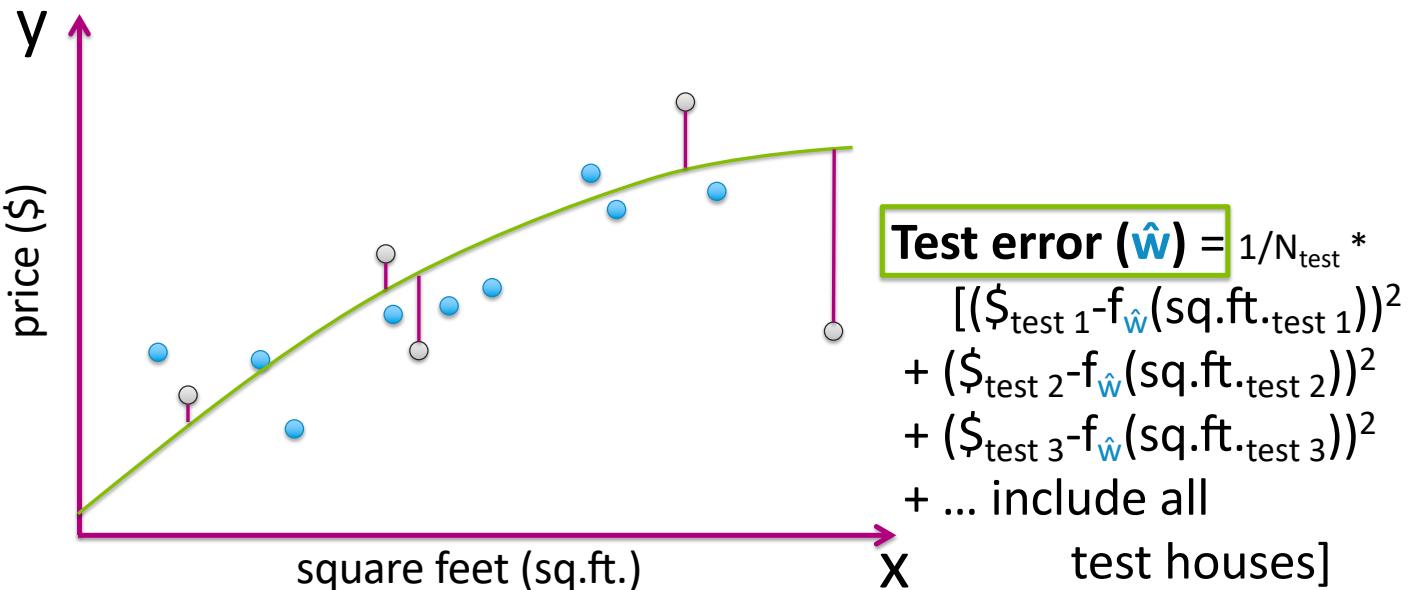
Example:

As before, fit quadratic to training data

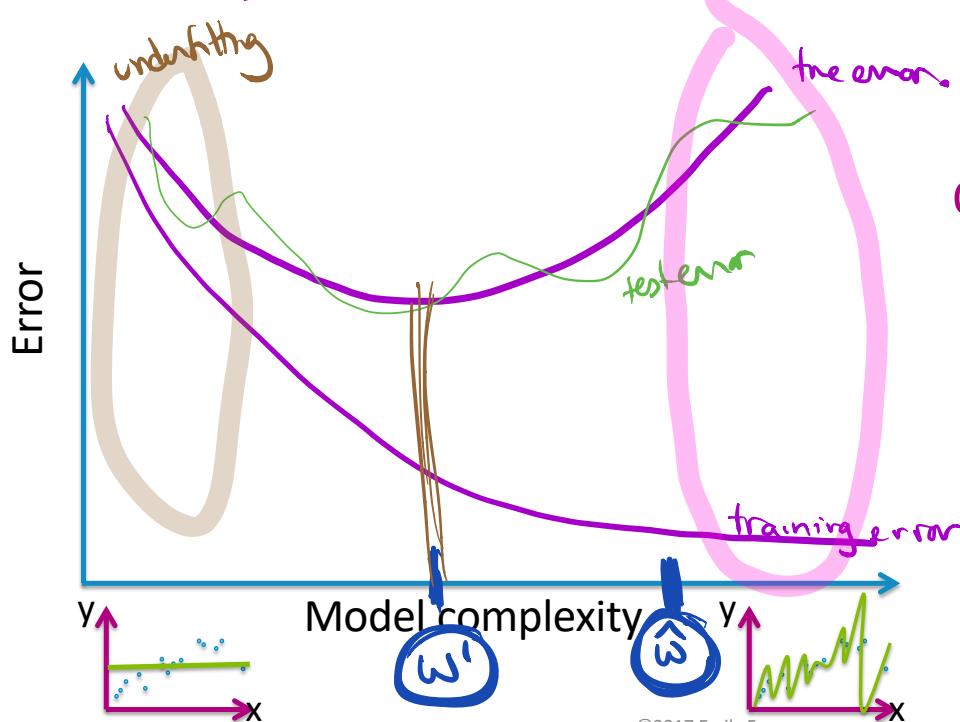


Assess performance using **test error**

As before, use **squared error loss** $(y - f_{\hat{w}}(x))^2$



Training, true, & test error vs. model complexity



Overfitting if:

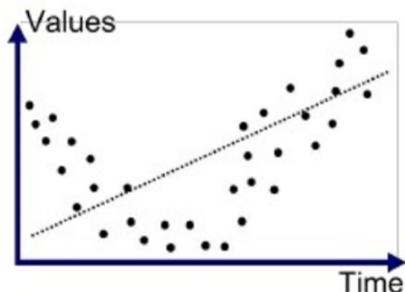
at \hat{w}

$y \exists w' \text{ s.t.}$

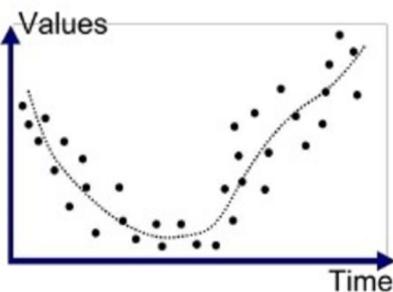
$\text{trainerr}(\hat{w}) < \text{trainerr}(w')$

$\text{testerr}(\hat{w}) > \text{testerr}(w')$

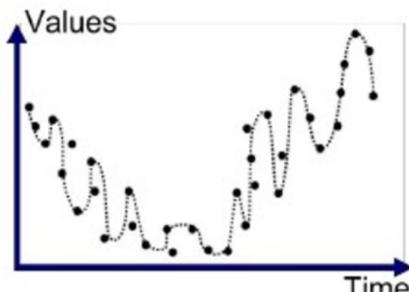
Underfitting vs Overfitting



Underfitted



Good Fit/R robust



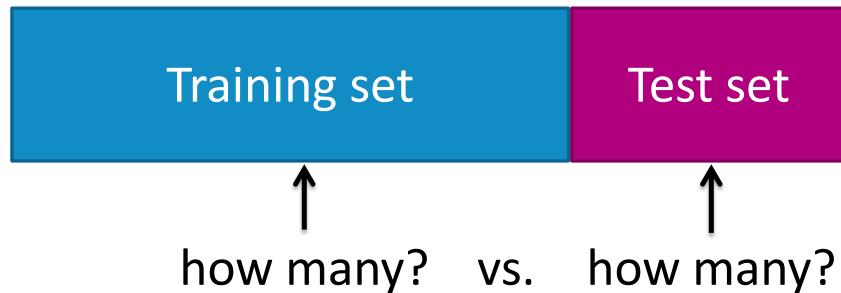
Overfitted



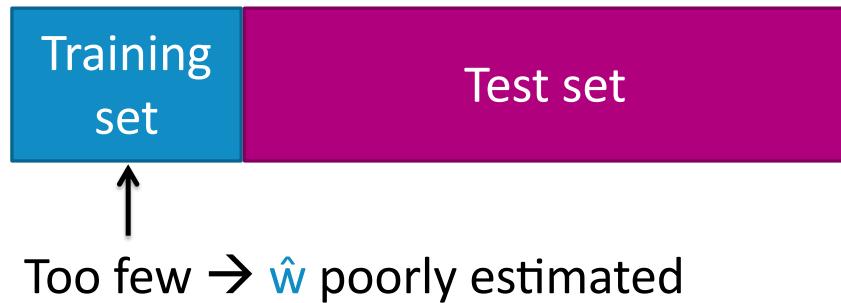
Courtesy Blog@AgoTrading101

Training/test split

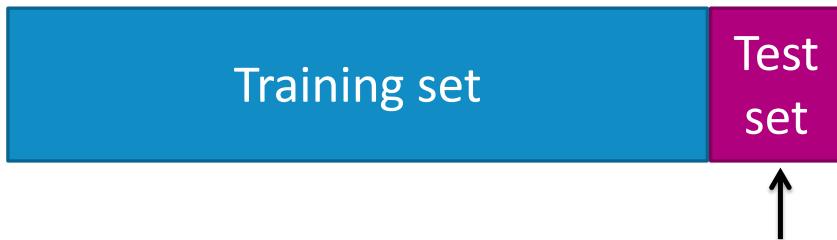
Training/test splits



Training/test splits

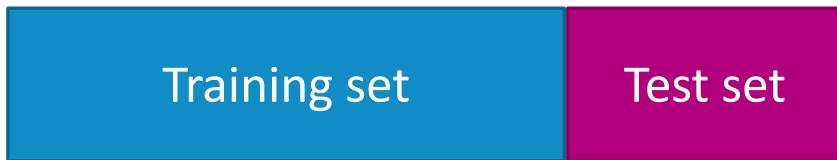


Training/test splits



Too few → test error bad approximation of
generalization error

Training/test splits



Typically, just enough test points to form a reasonable estimate of generalization error

If this leaves too few for training, other methods like cross validation (will see later...)

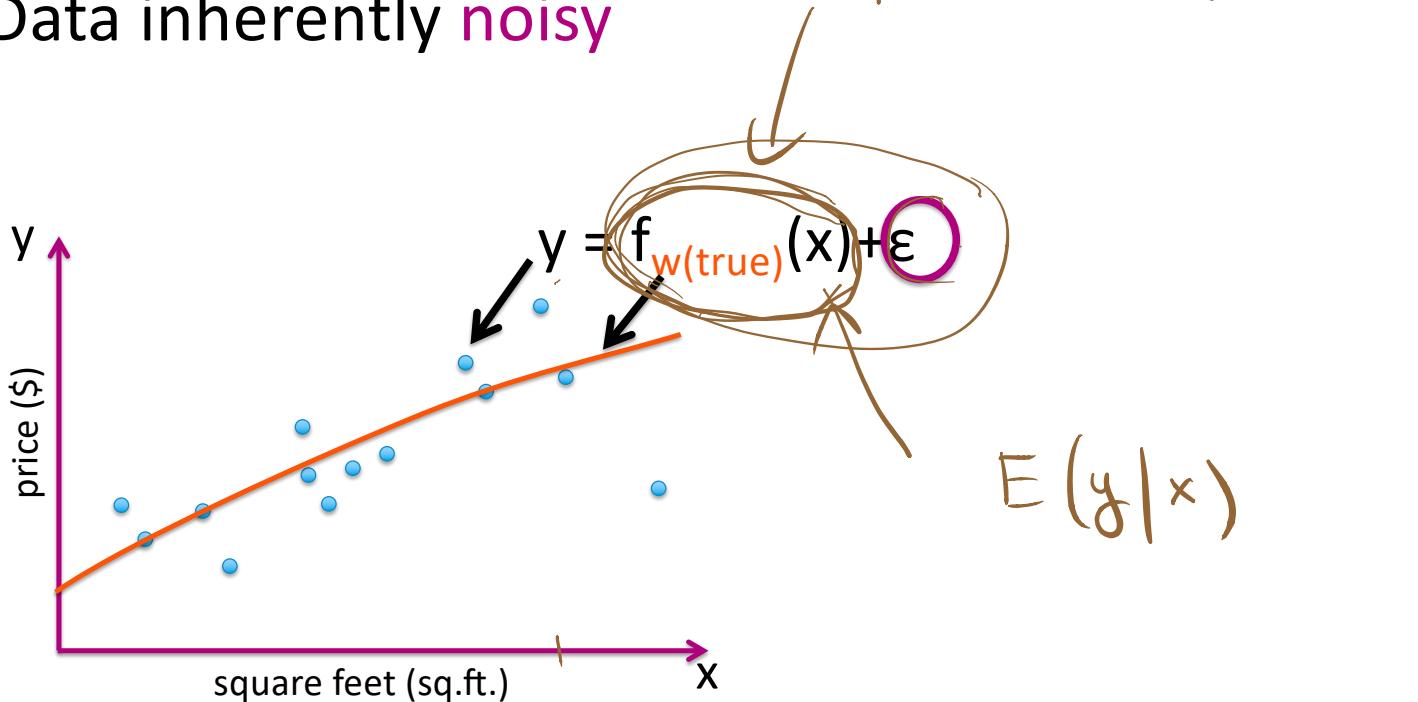
3 sources of error + the bias-variance tradeoff

3 sources of error

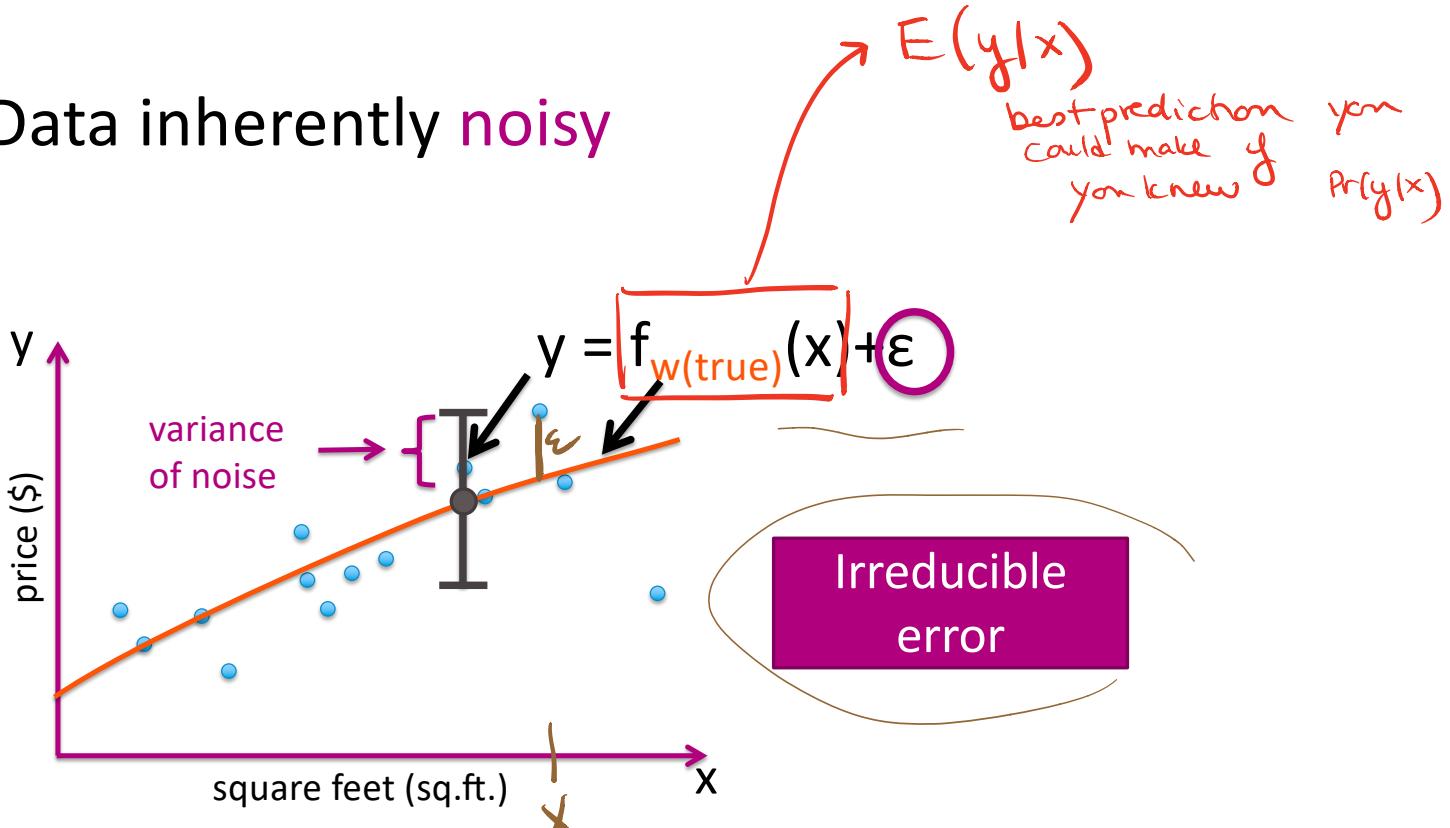
In forming predictions, there are 3 sources of error:

1. Noise
2. Bias
3. Variance

Data inherently **noisy**

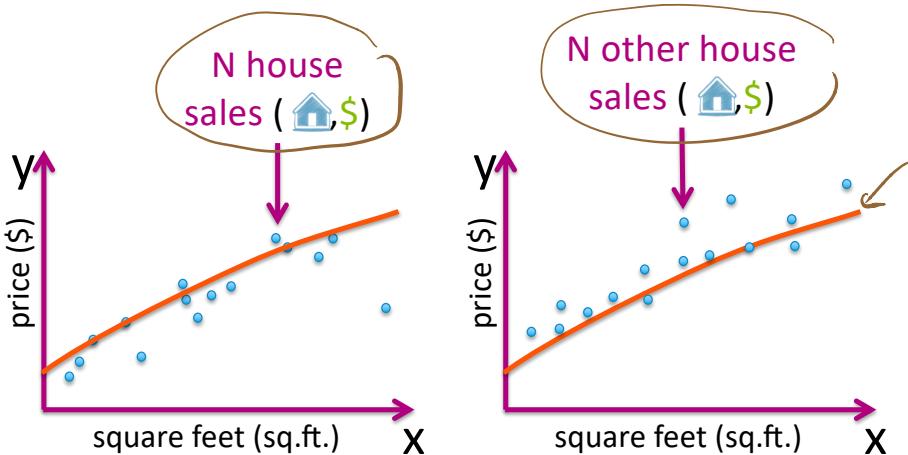


Data inherently noisy



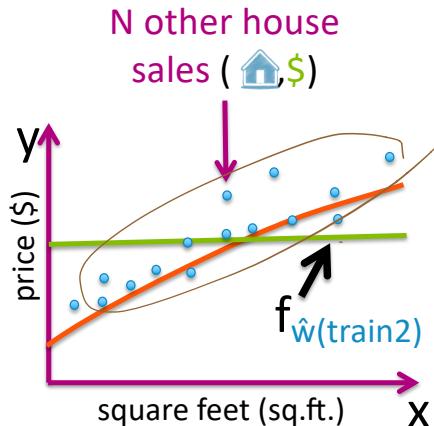
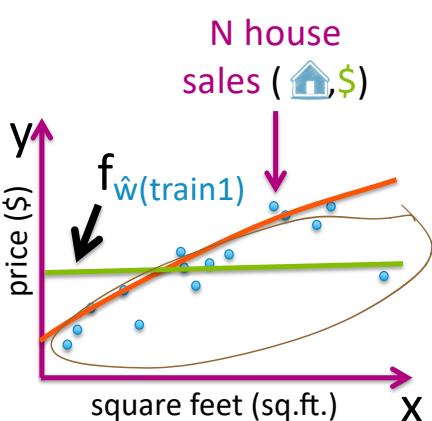
Bias contribution

Assume we fit a constant function



Bias contribution

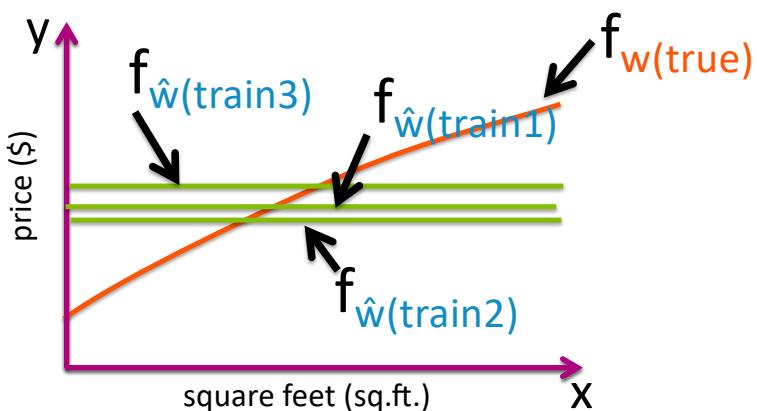
Assume we fit a constant function



$$f_{\bar{w}}(x) = E_{\text{train}} f[\hat{w}(\text{train})](x)$$

Bias contribution

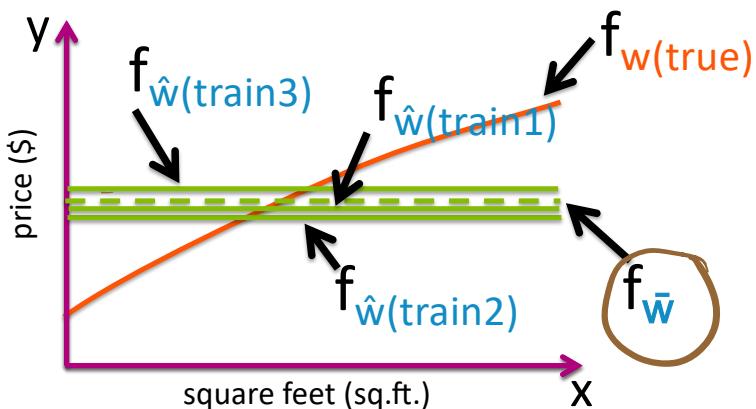
Over all possible size N training sets,
what do I expect my fit to be?



Bias contribution

$$f_{\bar{w}}(x) = E_{\text{train}} f[\hat{w}(\text{train})(x)]$$

Over all possible size N training sets,
what do I expect my fit to be?

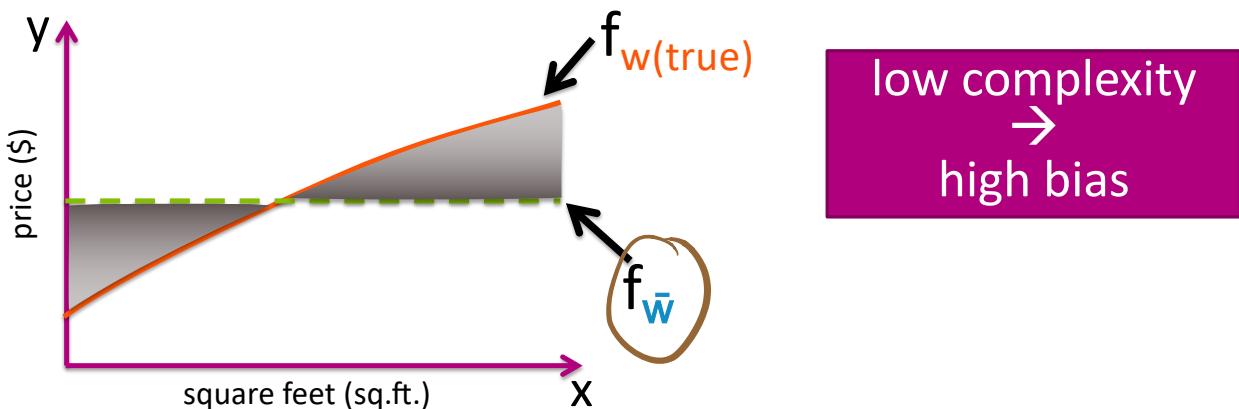


Bias contribution

$$\text{Bias}(x) = f_{w(\text{true})}(x) - f_{\bar{w}}(x)$$



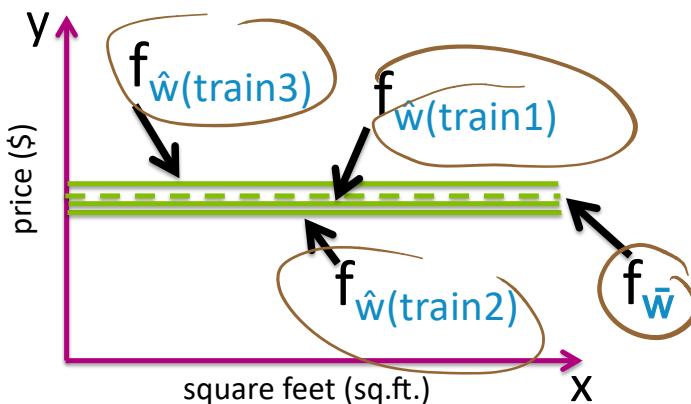
Is our approach flexible enough
to capture $f_{w(\text{true})}$?
If not, error in predictions.



Variance contribution

$$f_{\bar{w}}(x) = E_{\text{train}} f[\hat{w}(\text{train})(x)]$$

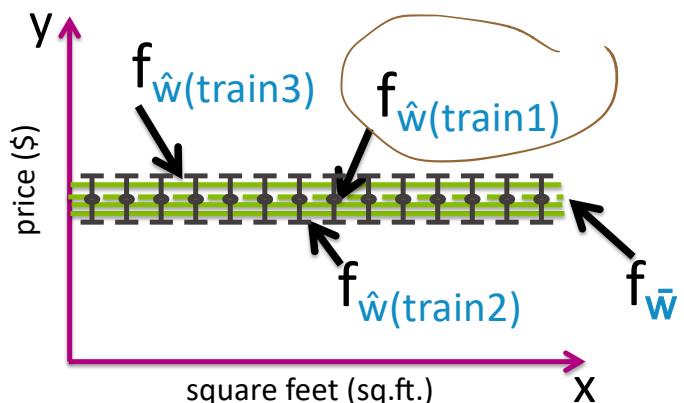
How much do specific fits vary from the expected fit?



$$f_{\bar{w}}(x) = E_{\text{train}} f[\hat{w}(\text{train})](x)$$

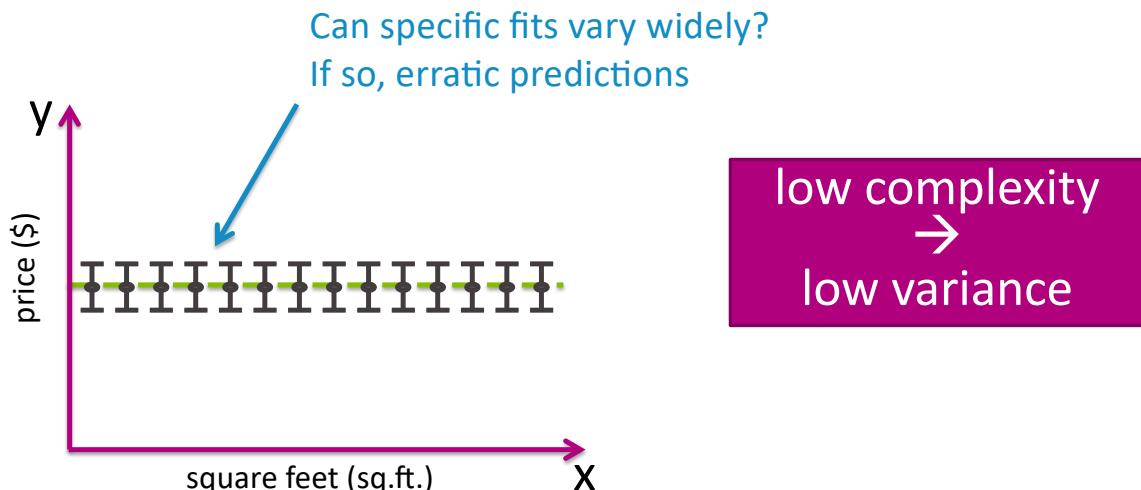
Variance contribution

How much do specific fits vary from the expected fit?



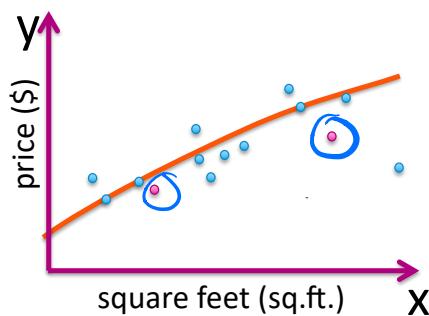
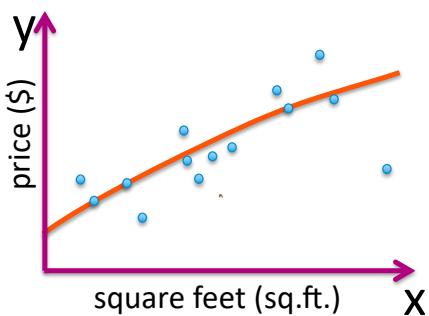
Variance contribution

How much do specific fits vary from the expected fit?



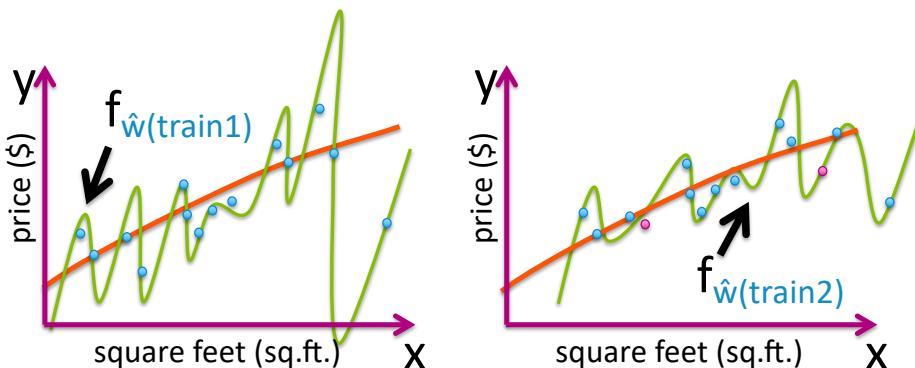
Variance of high-complexity models

Assume we fit a high-order polynomial



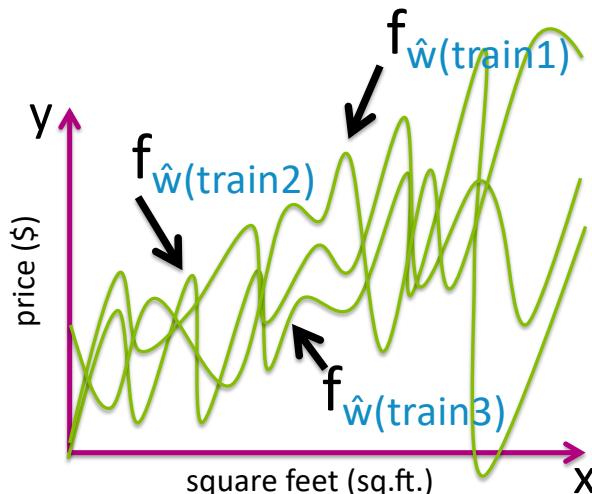
Variance of high-complexity models

Assume we fit a high-order polynomial



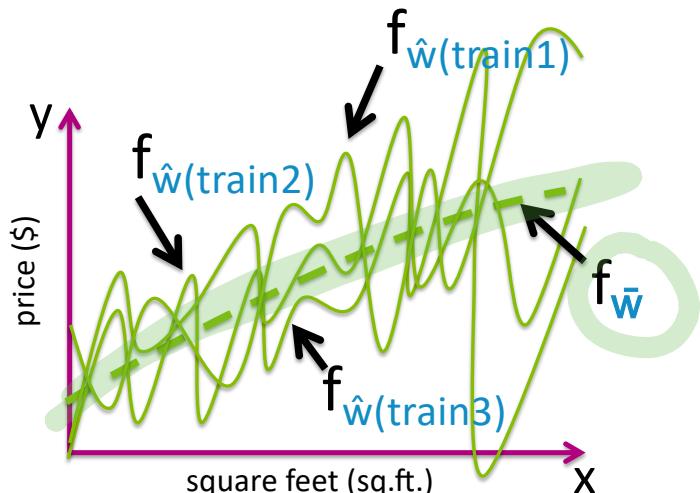
Variance of high-complexity models

Assume we fit a high-order polynomial

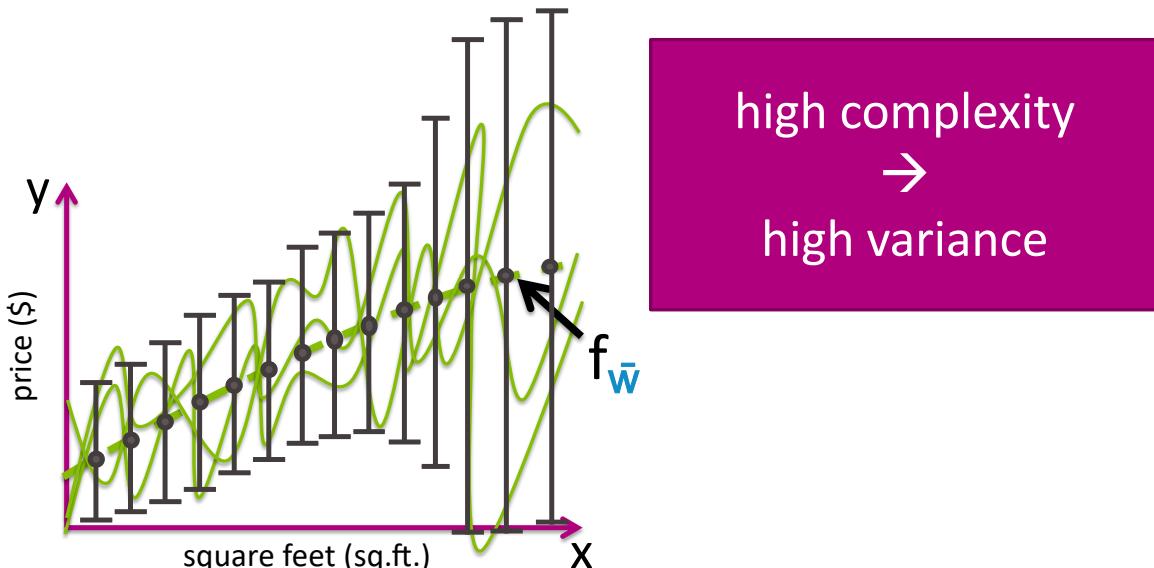


Variance of high-complexity models

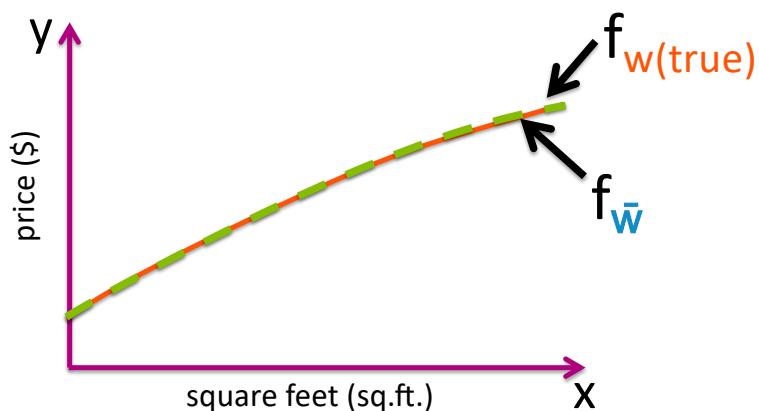
Assume we fit a high-order polynomial



Variance of high-complexity models

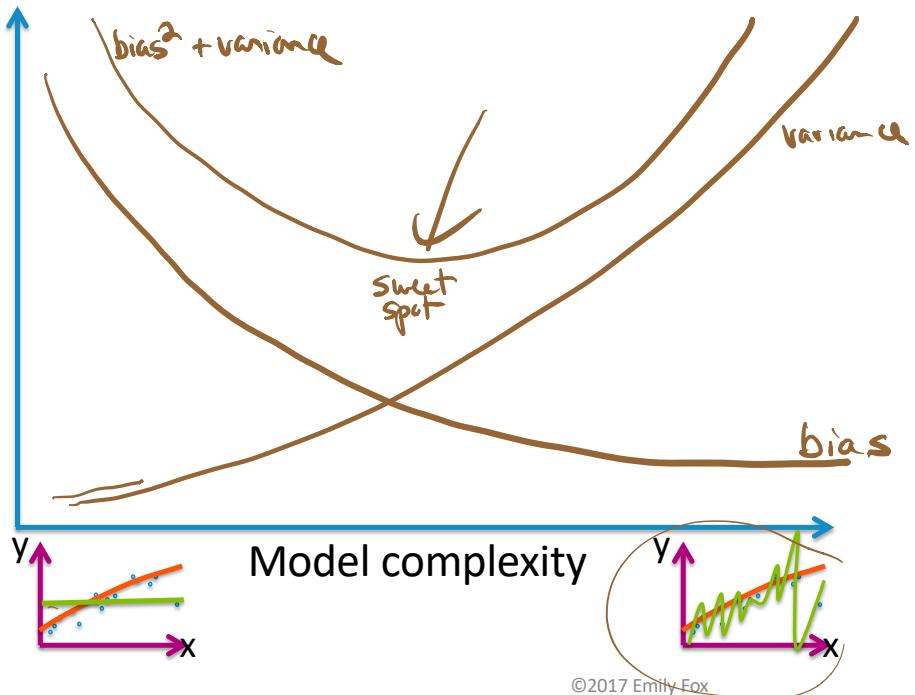


Bias of high-complexity models

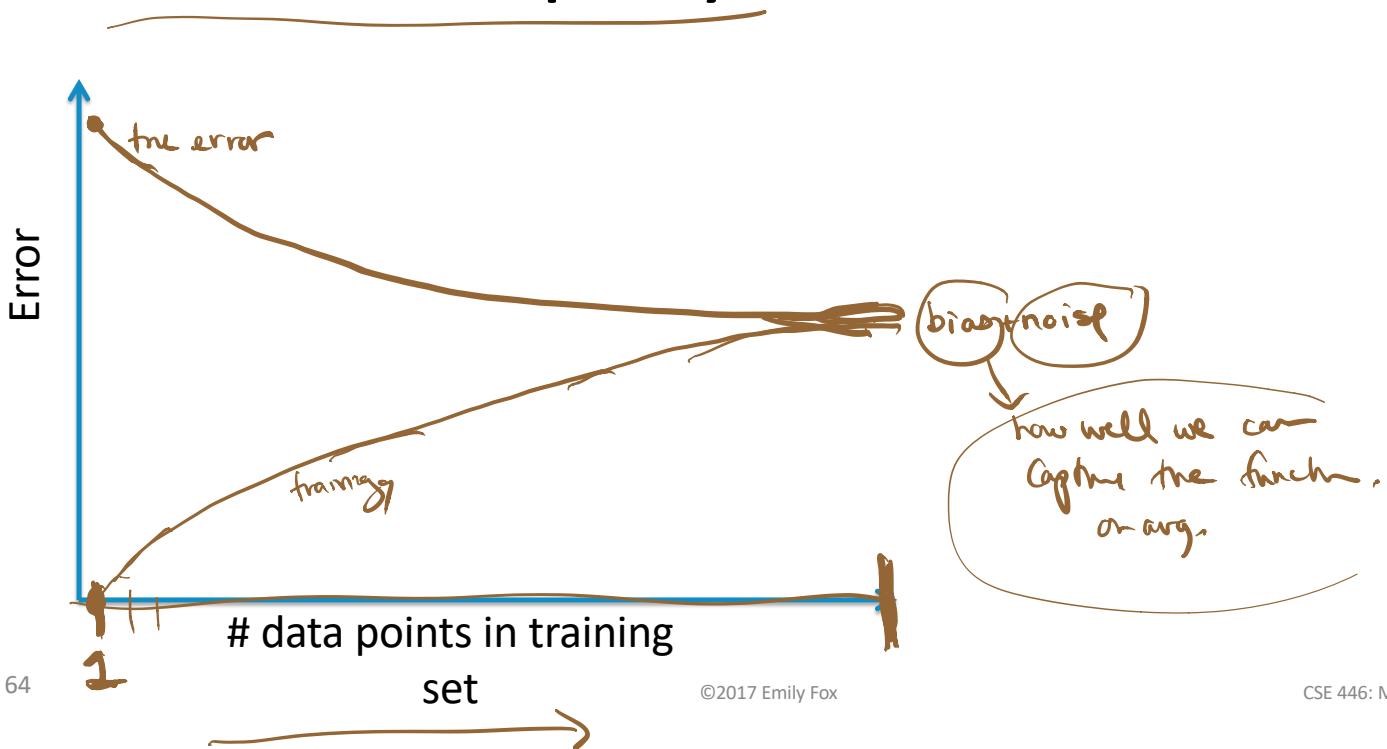


high complexity
→
low bias

Bias-variance tradeoff



True error and training error vs. amount of data for fixed model complexity

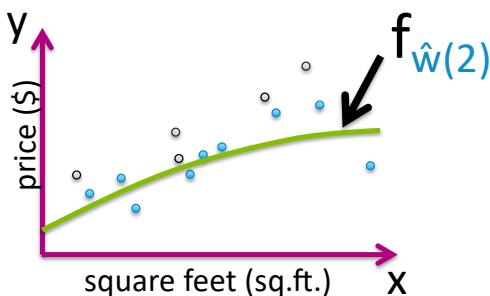
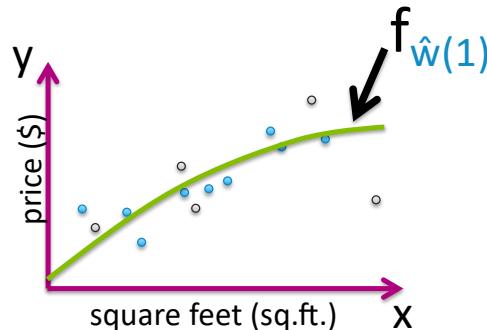


Formalize 3 sources of errors...

Accounting for training set randomness

Training set was just a random sample of n houses sold

What if n other houses had been sold and recorded?



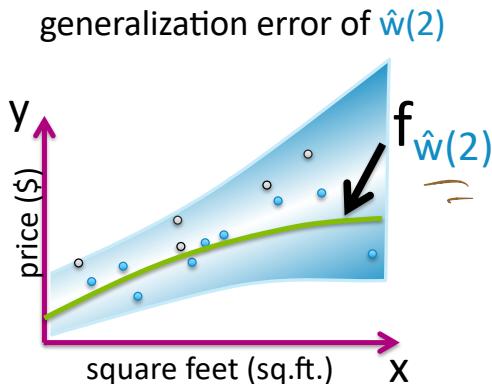
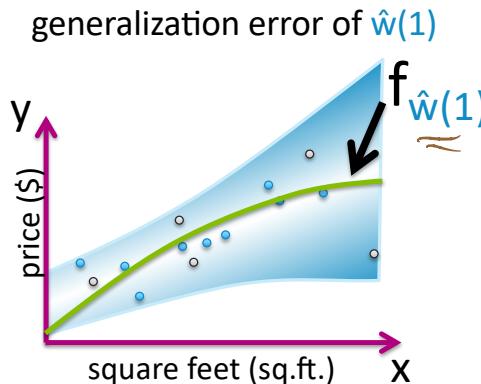
©2017 Emily Fox

CSE 446: Machine Learning

Accounting for training set randomness

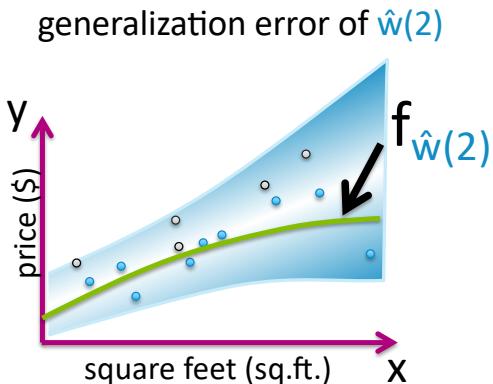
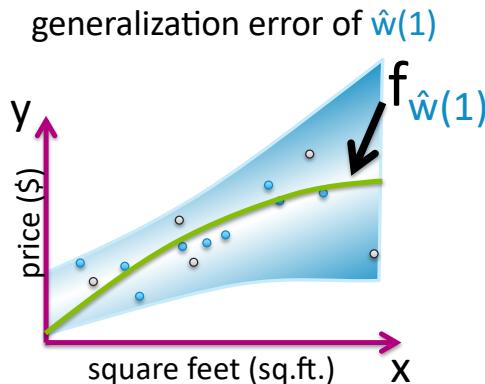
Training set was just a random sample of n houses sold

What if n other houses had been sold and recorded?



Accounting for training set randomness

Ideally, want performance averaged over all possible training sets of size n



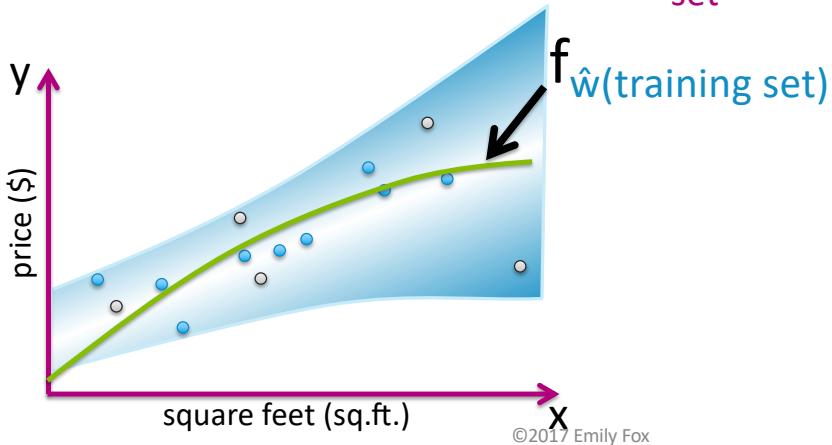
$$E_{\text{train}} = E_{x,y} \left[(y - f_{\hat{w}(\text{train})}(x))^2 \right]$$

Expected prediction error

$E_{\text{training set}}[\text{generalization error of } \hat{w}(\text{training set})]$

↑ averaging over all training sets
(weighted by how likely each is)

↑ parameters fit on
a specific training set



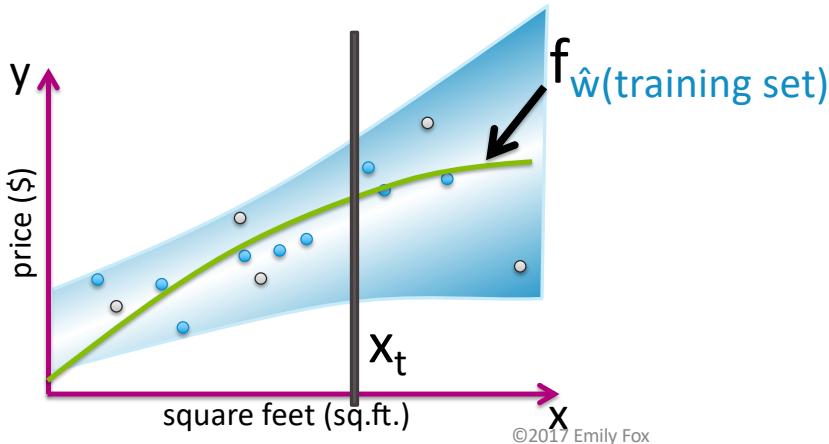
$$E_{\text{train}} = E_{x,y} \left[(y - f_{\hat{w}(\text{train})}(x))^2 \right]$$

= $E_x \left[E_{\text{train}, y|x} \left[(y - f_{\hat{w}(\text{train})}(x))^2 \right] \right]$

Prediction error at target input

Start by considering:

1. Loss at target x_t (e.g. 2640 sq.ft.)
2. Squared error loss $L(y, f_{\hat{w}}(x_t)) = (y - f_{\hat{w}}(x_t))^2$

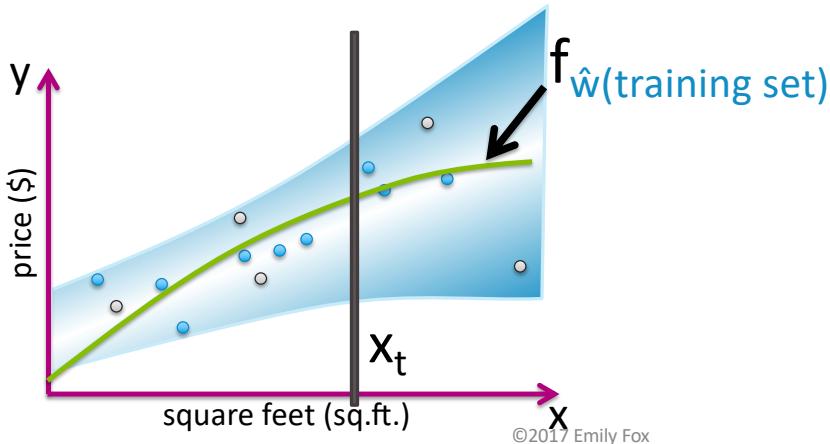


Sum of 3 sources of error

Average prediction error at x_t

$$= \sigma^2 + [\text{bias}(f_{\hat{w}}(x_t))]^2 + \text{var}(f_{\hat{w}}(x_t))$$

Variance of noise
Bias how well on avg model can fit $f_{w(\text{true})}$
Variance variability of fit

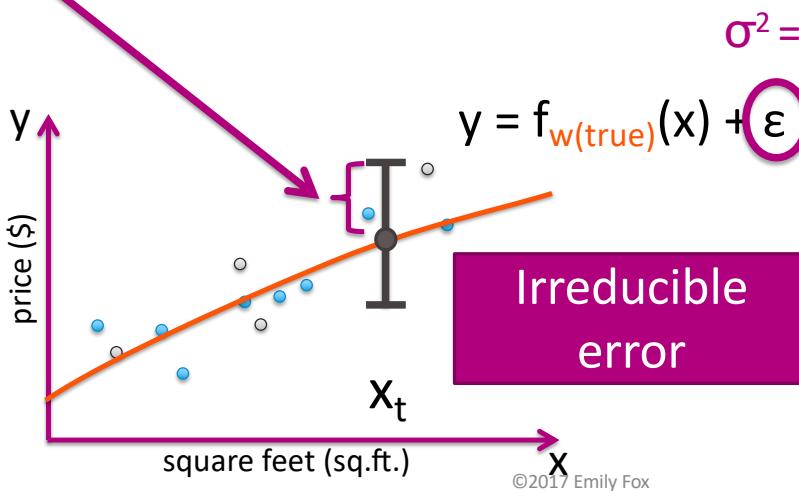


Error variance of the model

Variance of noise at x_t =
 $E_{y|x_t} [(y - f_{w(\text{true})}(x_t))^2] = \sigma^2$

Average prediction error at x_t

$$= \sigma^2 + [\text{bias}(f_{\hat{w}}(x_t))]^2 + \text{var}(f_{\hat{w}}(x_t))$$



σ^2 = "variance"

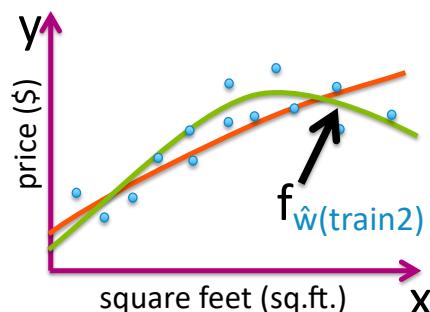
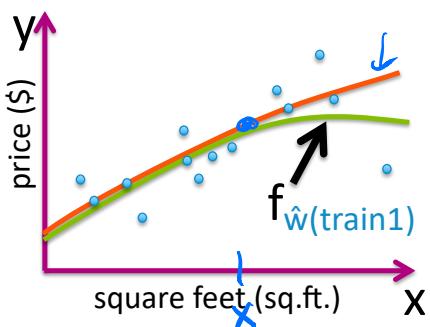
$$f_{w(\text{true})}(x) = E(y|x)$$

Bias of function estimator

Bias how well on avg model can fit $f_{\hat{w}(\text{true})}$
Variance variability of fit

Average prediction error at x_t

$$= \sigma^2 + [\text{bias}(f_{\hat{w}}(x_t))]^2 + \text{var}(f_{\hat{w}}(x_t))$$



Bias of function estimator

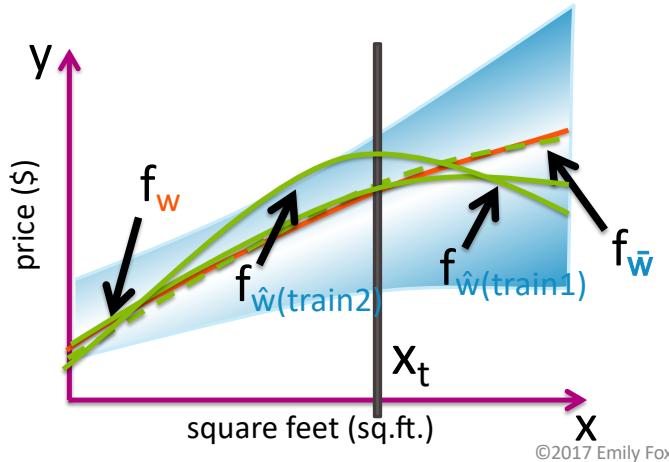
Average estimated function = $f_{\bar{w}}(x)$

True function = $f_w(x)$

\equiv

$$E_{\text{train}}[f_{\hat{w}(\text{train})}(x)]$$

over all training sets of size N

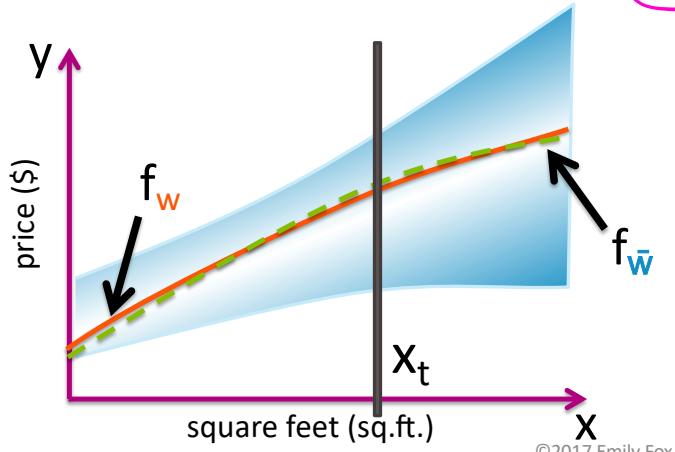


Bias of function estimator

Average estimated function = $f_{\bar{w}}(x)$

True function = $f_w(x)$

$$\text{bias}(f_{\hat{w}}(x_t)) = f_w(x_t) - f_{\bar{w}}(x_t)$$



Bias of function estimator

Average prediction error at x_t

$$= \sigma^2 + [\text{bias}(f_{\hat{w}}(x_t))]^2 + \text{var}(f_{\hat{w}}(x_t))$$

Average estimated function = $f_{\bar{w}}(x)$

True function = $f_w(x)$

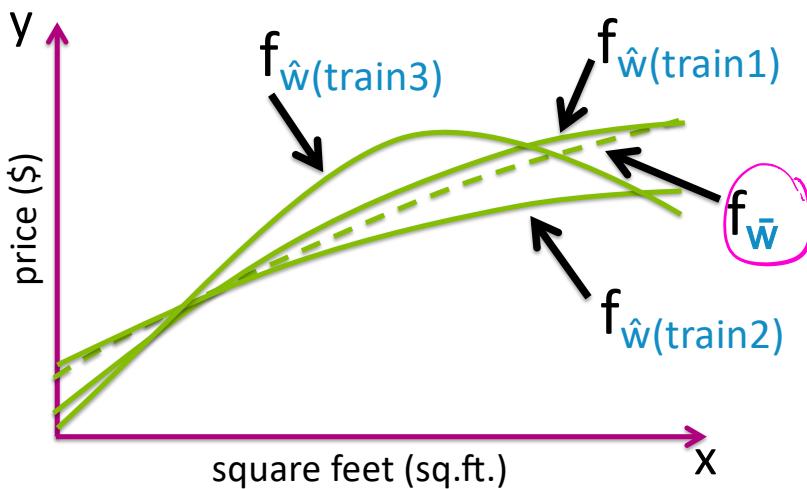
$$\text{bias}(f_{\hat{w}}(x_t)) = f_w(x_t) - f_{\bar{w}}(x_t)$$

Variance of function estimator

Bias how well on avg model can fit $f_{\hat{w}(\text{true})}$
Variance variability of fit

Average prediction error at x_t

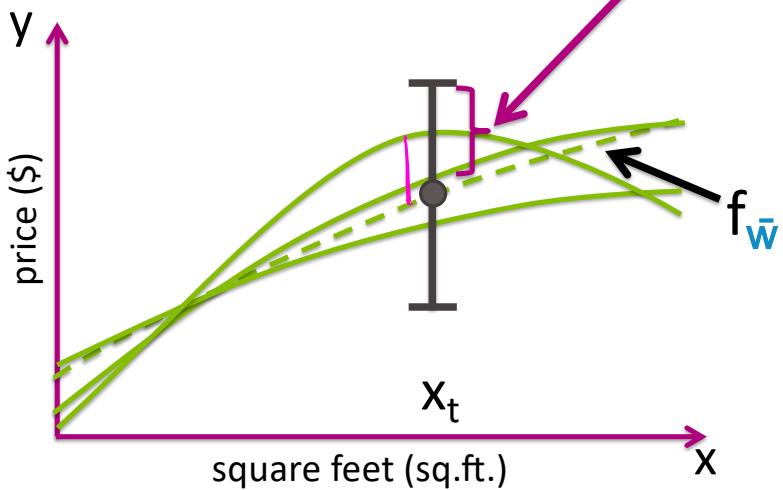
$$= \sigma^2 + [\text{bias}(f_{\hat{w}}(x_t))]^2 + \text{var}(f_{\hat{w}}(x_t))$$



Variance of function estimator

Average prediction error at x_t

$$= \sigma^2 + [\text{bias}(f_{\hat{w}}(x_t))]^2 + \text{var}(f_{\hat{w}}(x_t))$$

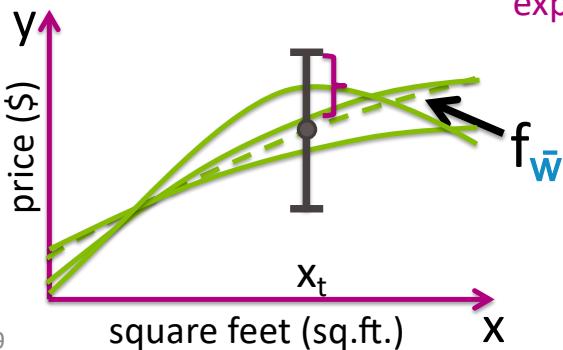


Variance of function estimator

$$\text{var}(f_{\hat{w}}(x_t)) = E_{\text{train}}[(f_{\hat{w}(\text{train})}(x_t) - f_{\bar{w}}(x_t))^2]$$

fit on a specific training dataset what I expect to learn over all training sets

over all training sets of size N deviation of specific fit from expected fit at x_t



$$f_{\bar{w}}(x_t) = \underset{\text{train}}{E} (f_{\hat{w}(\text{train})}(x_t))$$

Sum of 3 sources of error

Average prediction error at $x_t = E_{y|x}[(y - f_{\hat{w}(\text{train})}(x_t))^2]$

$$= \sigma^2 + [\text{bias}(f_{\hat{w}}(x_t))]^2 + \text{var}(f_{\hat{w}}(x_t))$$

$$\sigma^2 \text{ at } x_t = E_{y|x_t}[(y - f_{w(\text{true})}(x_t))^2]$$

$$\text{Bias}(x_t) = f_{w(\text{true})}(x_t) - f_{\bar{w}}(x_t)$$

$$\text{Variance}(x_t) = E_{\text{train}}[(f_{\hat{w}(\text{train})}(x_t) - f_{\bar{w}}(x_t))^2]$$

$$f_{\bar{w}}(x_t) = E_{\text{train}} f[\hat{w}(\text{train})(x_t)]$$

Why 3 sources of error?
A formal derivation

Deriving expected prediction error

Expected prediction error

$$= E_{\text{train}} [\text{generalization error of } \hat{w}(\text{train})]$$

$$= E_{\text{train}} [E_{x,y} [L(y, f_{\hat{w}(\text{train})}(x))]]$$

1. Look at specific x_t
2. Consider $L(y, f_{\hat{w}}(x)) = (y - f_{\hat{w}}(x))^2$

Expected prediction error at x_t

$$= E_{\text{train}, y|x_t} [(y - f_{\hat{w}(\text{train})}(x_t))^2]$$



Deriving expected prediction error

Expected prediction error at x_t

$$= E_{\text{train}, y|x_t} [(y - f_{\hat{w}(\text{train})}(x_t))^2]$$

$$= E_{\text{train}, y|x_t} [((y - f_{w(\text{true})}(x_t)) + (f_{w(\text{true})}(x_t) - f_{\hat{w}(\text{train})}(x_t)))^2] = E(C^2) + E(CD) + E(D^2)$$

C D

$$E(C^2) = \boxed{E_{y|x_t} [(y - f_{w(\text{true})}(x_t))^2]} = \sigma^2 \quad \text{noise}$$

$$E(CD) = E_{\text{train}} [(f_{w(\text{true})}(x_t) - f_{\hat{w}(\text{train})}(x_t)) E_{y|x_t} (y - f_{w(\text{true})}(x_t))] = 0 \quad \text{since } f_{w(\text{true})}(x_t) = E(y|x_t)$$

$$E(D^2) = E_{\text{train}} \left[(f_{w(\text{true})}(x_t) - f_{\hat{w}(\text{train})}(x_t))^2 \right] \rightarrow \text{continued on next page.}$$

$\underbrace{(f_{w(\text{true})}(x_t) - f_{\hat{w}(\text{train})}(x_t))^2}_{\text{MSE}[f_{\hat{w}(\text{train})}(x_t)]}$

Equating MSE with bias and variance

$$f_{\bar{w}}(x) = E_{\text{train}} f[\hat{w}(\text{train})(x)]$$

$$\text{MSE}[f_{\hat{w}(\text{train})}(x_t)]$$

$$= E_{\text{train}} [(f_{w(\text{true})}(x_t) - f_{\hat{w}(\text{train})}(x_t))^2]$$

$$= E_{\text{train}} [\underbrace{(f_{w(\text{true})}(x_t) - f_{\bar{w}}(x_t))}_{A} + \underbrace{(f_{\bar{w}}(x_t) - f_{\hat{w}(\text{train})}(x_t))}_{B}]^2$$

$$E_{\text{train}}(A^2) = E_{\text{train}}[(f_{w(\text{true})}(x_t) - f_{\bar{w}}(x_t))^2] = \boxed{(f_{w(\text{true})}(x_t) - f_{\bar{w}}(x_t))^2 = (\text{bias}(x_t))^2}$$

$$E_{\text{train}}(AB) = (f_{w(\text{true})}(x_t) - f_{\bar{w}}(x_t)) E_{\text{train}}(f_{\bar{w}}(x_t) - f_{\hat{w}(\text{train})}(x_t))$$

$$= 0 \quad \text{since } E_{\text{train}}(f_{\hat{w}(\text{train})}(x_t)) = f_{\bar{w}}(x_t)$$

$$\begin{aligned} & E((A+B)^2) \\ &= E(A^2) + 2E(AB) + E(B^2) \end{aligned}$$

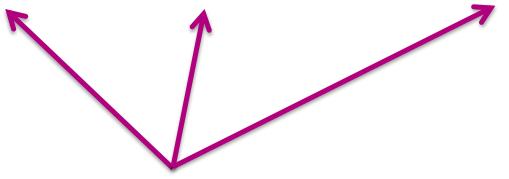
$$E_{\text{train}}(B^2) = \boxed{E_{\text{train}}[(f_{\hat{w}(\text{train})}(x_t) - f_{\bar{w}}(x_t))^2]} = \text{variance of } f_{\hat{w}}(x_t)$$

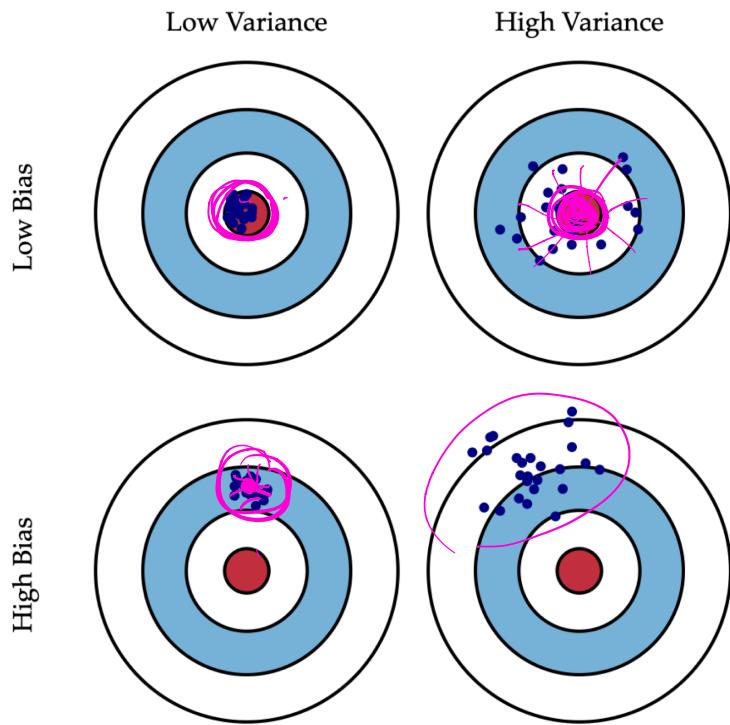
Putting it all together

Expected prediction error at x_t

$$= \sigma^2 + \text{MSE}[f_{\hat{w}}(x_t)]$$

$$= \sigma^2 + [\text{bias}(f_{\hat{w}}(x_t))]^2 + \text{var}(f_{\hat{w}}(x_t))$$





Courtesy of Scott Fortmann-Roe

Summary of assessing performance

What you can do now...

- Describe what a loss function is and give examples
- Contrast training, generalization, and test error
- Compute training and test error given a loss function
- Discuss issue of assessing performance on training set
- Describe tradeoffs in forming training/test splits
- Understand the 3 sources of avg. prediction error
 - Irreducible error, bias, and variance