

Assessing Performance

Measuring loss

Loss function

$$L(y, \hat{f}_w(x))$$

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

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

Examples: (assuming loss for underpredicting = overpredicting)

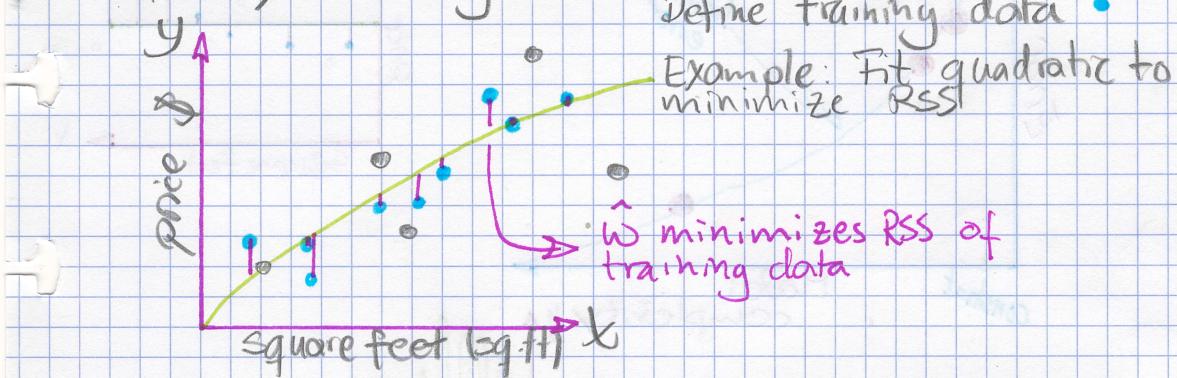
Absolute error: $L(y, \hat{f}_w(x)) = |y - \hat{f}_w(x)|$

Squared error: $L(y, \hat{f}_w(x)) = (y - \hat{f}_w(x))^2$

Training error: assessing loss on the training set

Assessing the loss

Part 1: Training error



Compute training error

1. Define a loss function $L(y, \hat{f}_w(x))$
- Eg squared error, absolute error, ...

2. Training error
= avg loss on houses in training set

$$= \frac{1}{N} \sum_{i=1}^N L(y_i, \hat{f}_w(x_i))$$

fit using training data

Example:

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

Training error (\hat{w}) =

$$= \frac{1}{N} \cdot [(\$_{train1} - \hat{f}_{\hat{w}}(\text{sq ft train1}))^2 +$$

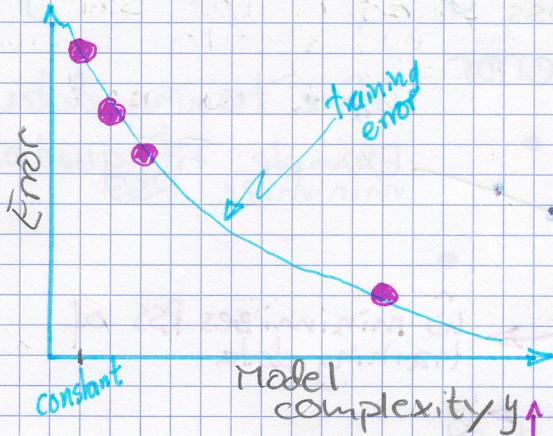
$$+ (\$_{train2} - \hat{f}_{\hat{w}}(\text{sq ft train2}))^2 +$$

+ ... + include all training houses]

$$= \frac{1}{N} \sum_{i=1}^n (y_i - \hat{f}_{\hat{w}}(x_i))^2$$

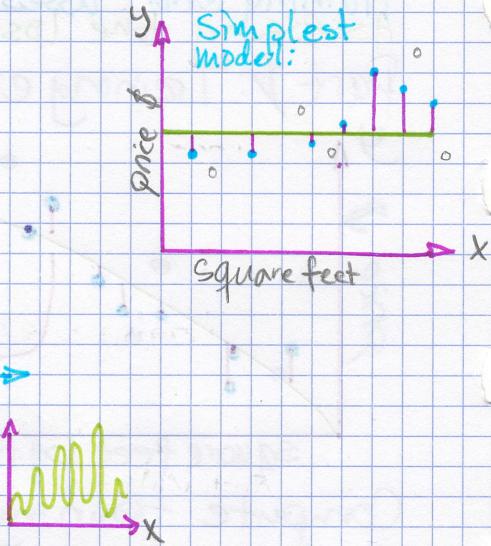
$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^n (y_i - \hat{f}_{\hat{w}}(x_i))^2}$$

{Root mean square error}

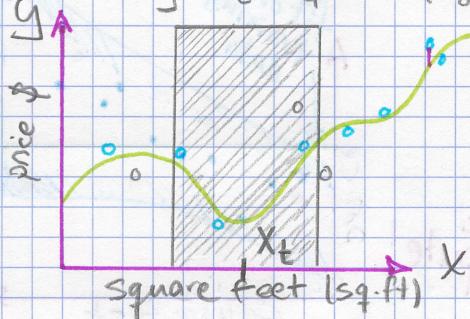


Is training error a good measure of predictive performance?

How do we expect to perform on a new house?



Is there something particularly bad about having x_t square feet???



Issue: Training error is overly optimistic because it was fit to training data

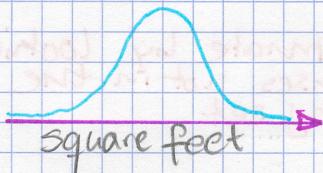
{ Small training error
⇒ good predictions
unless training data includes everything you might ever see

Assessing the loss

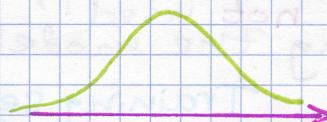
Part 2: Generalization (true) error

Distribution over houses

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



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



For fixed #sq.ft

Generalization error definition

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

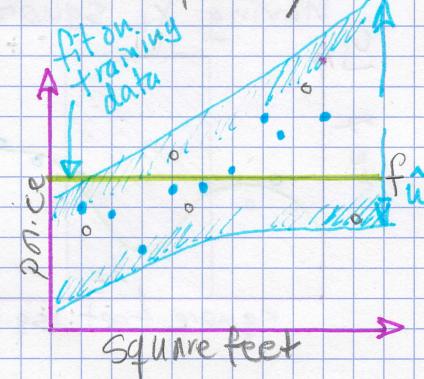
Average over all possible (x, y) pairs weighted by how likely each is

Finally:

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

fitting using training data

Generalization error vs. model complexity



Assessing the loss

Part 3: Test error

Approximating generalization error
Wanted estimate of loss over all possible (\hat{m} , \hat{s}) pairs

Forming a test set
Hold out some (\hat{m} , \hat{s})
that are not used
for fitting the model

↑ { Approximate by looking
at houses not in the
training set }



Compute test error

Test error

= avg. loss on houses in test set

$$= \frac{1}{N_{\text{test}}} \sum_{j \in \text{test set}} L(y_j, f_w^*(x_j))$$

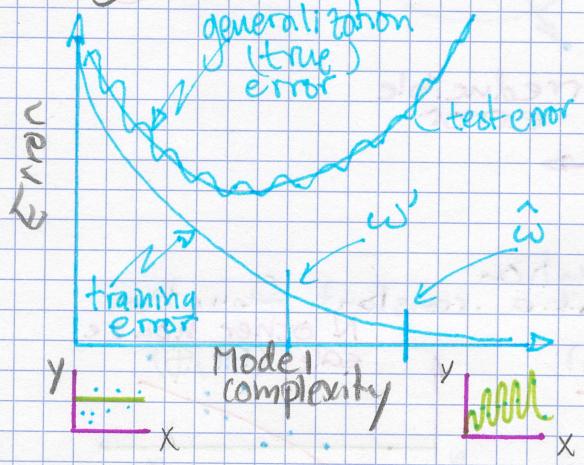
test points

test set

↑ fit using training data

$$\text{Test error } (\hat{\omega}) = 1/N * \left[(\$_{\text{test1}} - \hat{f}_{\omega}(\text{sq.ft. test1}))^2 + (\$_{\text{test2}} - \hat{f}_{\omega}(\text{sq.ft. test2}))^2 + (\$_{\text{test3}} - \hat{f}_{\omega}(\text{sq.ft. test3}))^2 + \dots \text{ include all test houses} \right]$$

Training, true, & test error vs. Model complexity



Overfitting if:

if there exists a model with estimated params $\hat{\omega}'$ such that

① training error ($\hat{\omega}$) $<$ training error ($\hat{\omega}'$)

② true error ($\hat{\omega}'$) $>$ true error ($\hat{\omega}$)

Training / test split

Training set	Test set
--------------	----------

Too few $\rightarrow \hat{\omega}$ poorly estimated

Training set	Test set
--------------	----------

Too few \rightarrow test error bad approximation of generalization error

Training set	Test set
--------------	----------

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

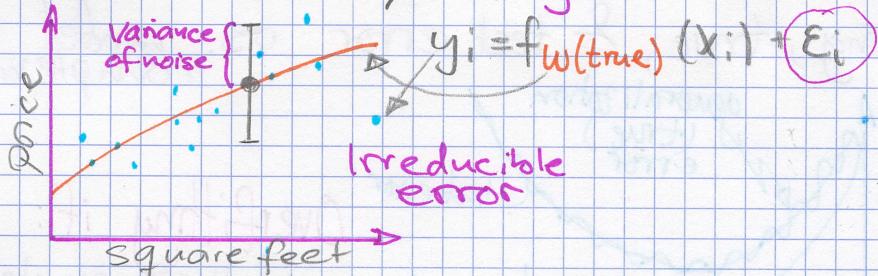
Irreducible error and bias

3 sources of error

In forming predictions, there are 3 sources of error

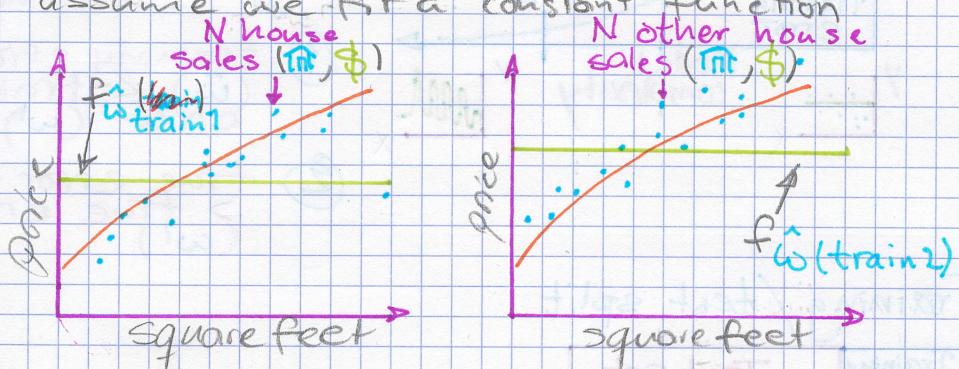
1. Noise
2. Bias
3. Variance

Data inherently noisy

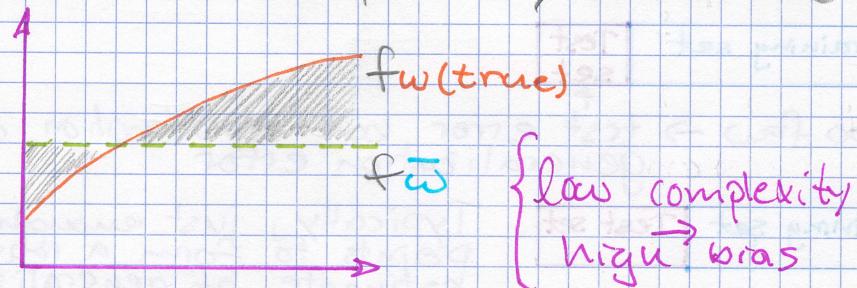


Bias contribution

assume we fit a constant function



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

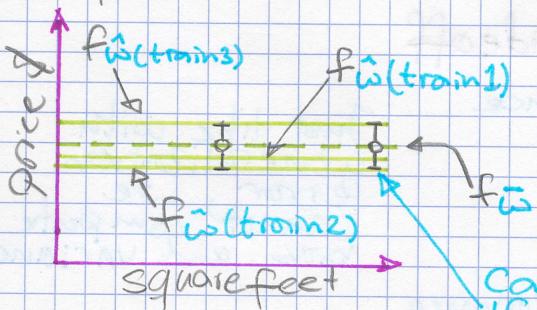


$$\text{Bias}(x) = f_{\omega(\text{true})}(x) - \bar{f}_{\omega}$$

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

Variance contribution

How much do specific fits vary from the expected fit



low complexity
→ low variance

Can specific fits vary widely?
If so, erratic predictions

Variance of high-complexity models
assume we fit a high order polynomial



high complexity

high variance

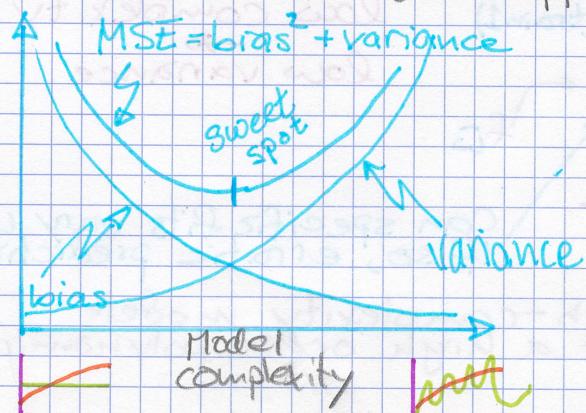


Bias of high-complexity models



High complexity
low bias

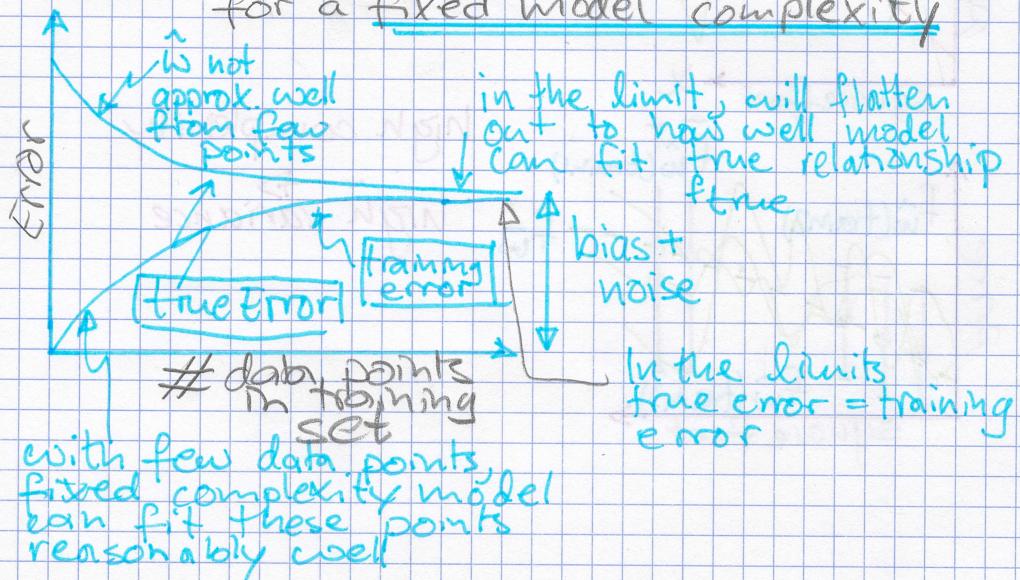
Bias - variance tradeoff



Just like with generalization error we cannot compute bias and variance

ERROR VS. AMOUNT OF DATA

for a fixed model complexity



The regression/ML workflow

1. Model selection

Often, need to choose tuning parameters λ controlling model complexity (e.g. degree of polynomial)

2. Model assessment

Having selected a model, assess the generalization error

Hypothetical implementation



1. Model selection

For each considered model complexity λ :

- i. Estimate parameters \hat{w}_λ on training data
- ii. Assess performance of \hat{w}_λ on test data
- iii. Choose λ^* to be λ with lowest test error

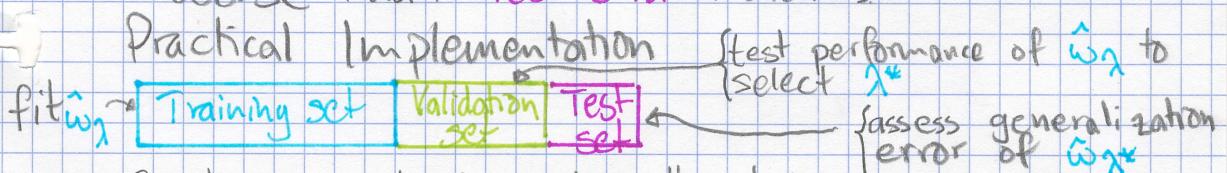
2. Model assessment

Overly optimistic!
Compute test error of \hat{w}_{λ^*} (fitted model for selected complexity λ^*) to approx. generalization error.

Issue: Just like fitting \hat{w} and assessing its performance both on training data

- λ^* was selected to minimize test error (i.e., λ^* was fit on test data)
- if test data is not representative of the whole world, then \hat{w}_{λ^*} will typically perform worse than test error indicates

Practical Implementation



Solution: Create two "test" sets!

1. Select λ^* such that \hat{w}_{λ^*} minimizes error on validation set

2. Approximate generalization error of \hat{w}_g^*
using test set

Typical splits

	Training set	Validation set	Test set
1)	60%	10%	10%
2)	50%	25%	25%