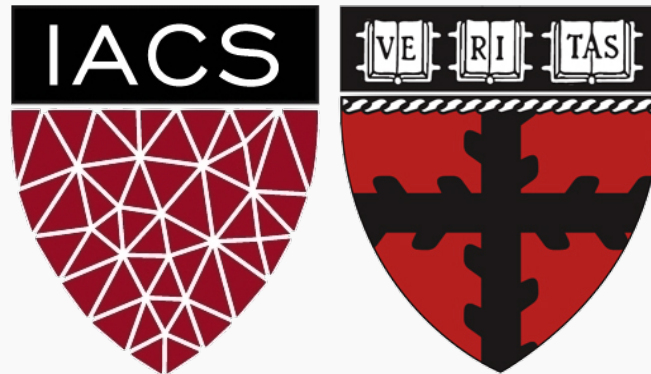# Generalization Error and Bias Variance Tradeoff

## CS109A Introduction to Data Science
Pavlos Protopapas, Natesh Pillai
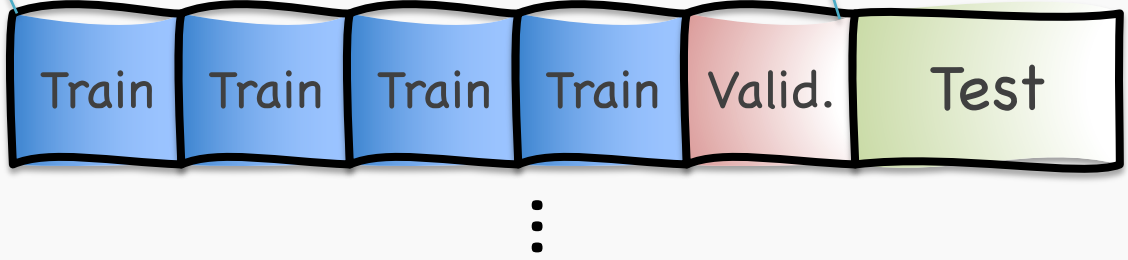
# Outline

- Q&A from lecture 5:

    - Train/Validation/Test

    - Scaling

- Generalization Error, Bias Variance Tradeoff

- Regularization

    o Lasso and Ridge

In the beginning, we always separate a portion of the data from the main dataset, which we never touch until the very end when we want to evaluate the performance of the final model. Normally, this is called train + test split. *

Sometimes we can further split train data into train + validation, essentially ending up with train + validation (both used to find the best model) **+ test** (which we still don't use until the very end).

And then, we sometimes also use cross-validation, which has nothing to do with either test or validation splits? Because cross-validation uses the train data to split it into k buckets.

* sometimes they (not CS109A) also call this train + validation split, while meaning train + test

When you realize k-Fold Cross Validation can only validate your hyperparameters, not yourself..

# Previously on CS109A

# Model Selection

1. Model selection as a way to avoid overfitting
2. Validation set to select the best model
3. Cross validation to avoid overfitting to the validation set

Ways of model selection:

- Exhaustive search
- Greedy algorithms
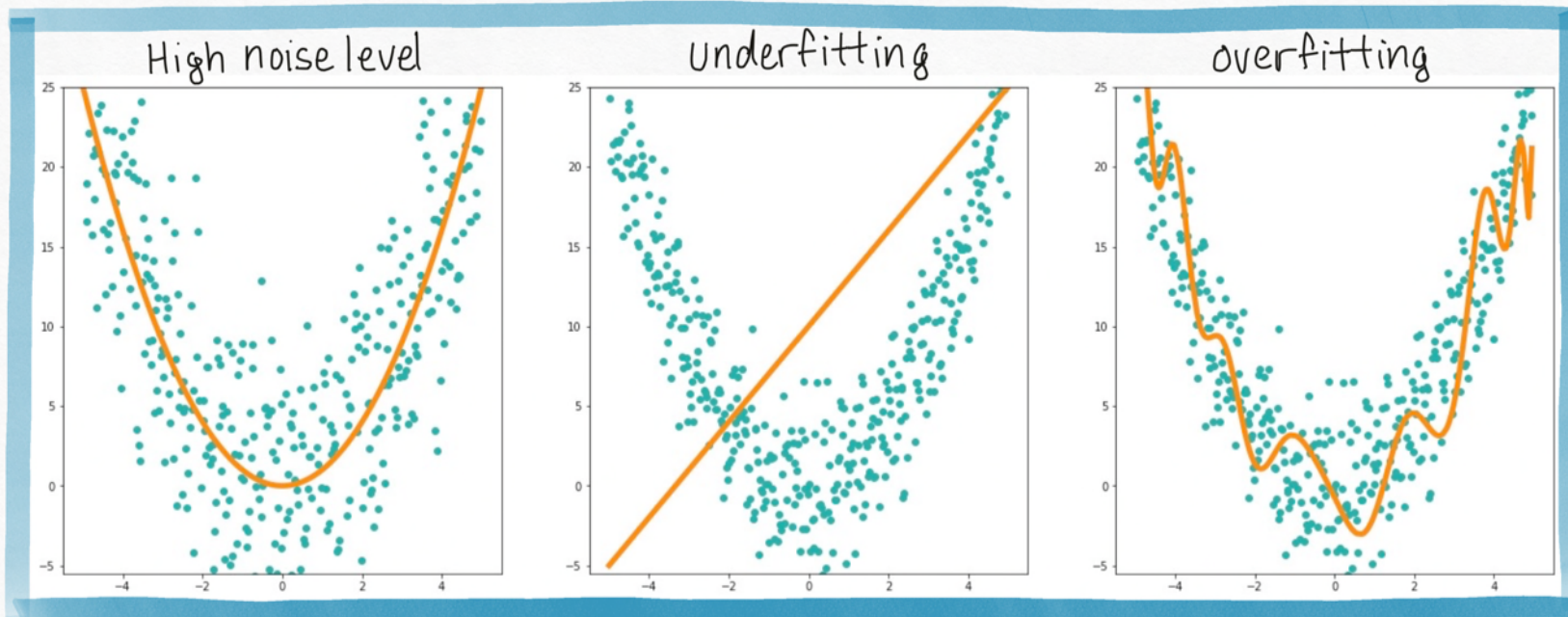- Fine tuning hyper-parameters
- **Regularization**

# Outline

- Q&A from lecture 5:

  - Train/Validation/Test

  - Scaling

- **Generalization Error, Bias Variance Tradeoff**

- Regularization

  - Lasso and Ridge

# Test Error and Generalization

We know to evaluate models on both train and test data because models can do well on training data but do poorly on new data.

When models do well on new data is called generalization.

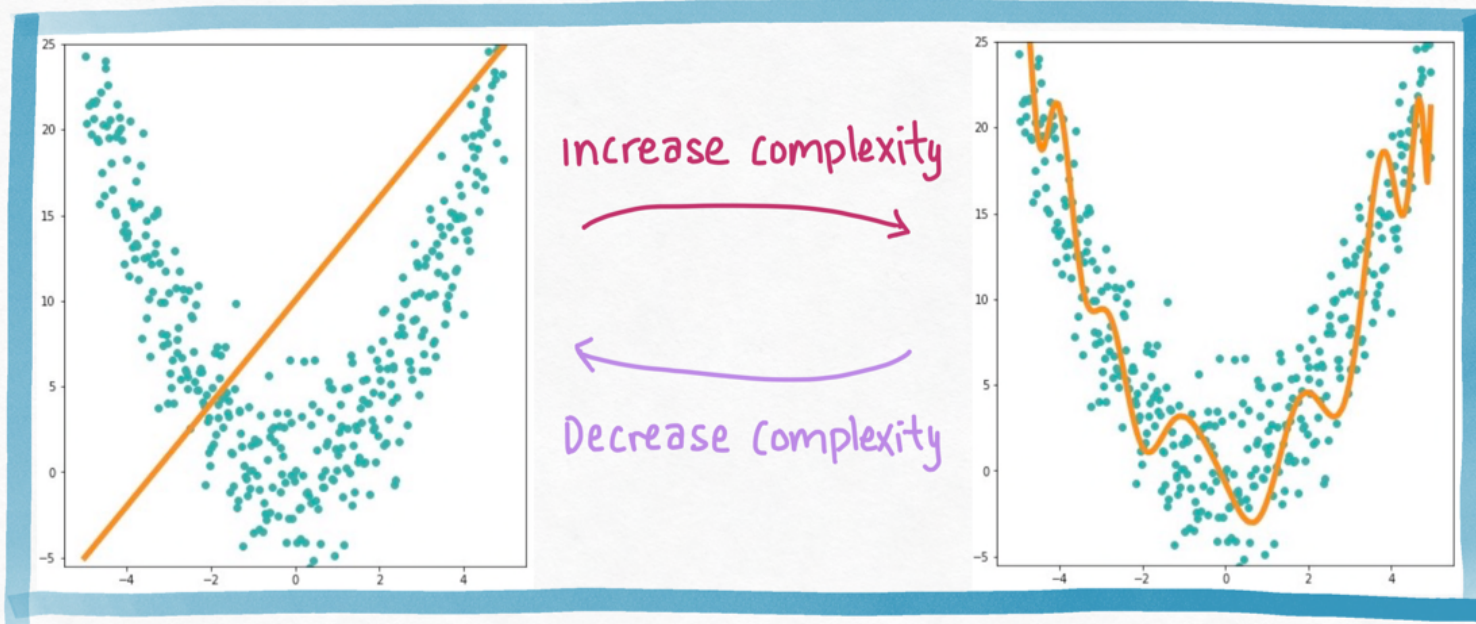There are at least three ways a model can have a high test error.

# Irreducible and Reducible Errors

We distinguished the contributions of noise to the generalization error:

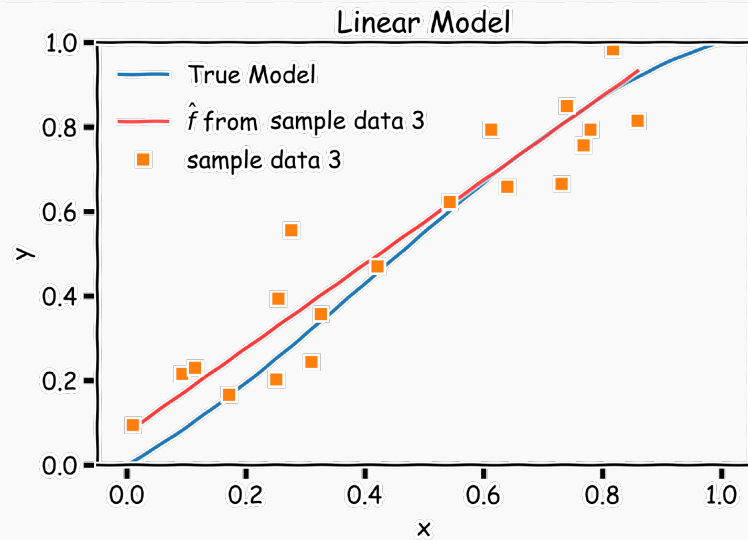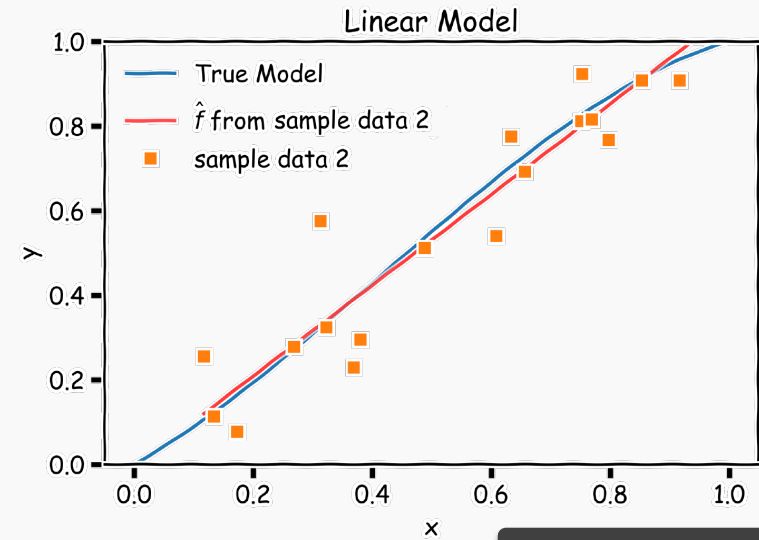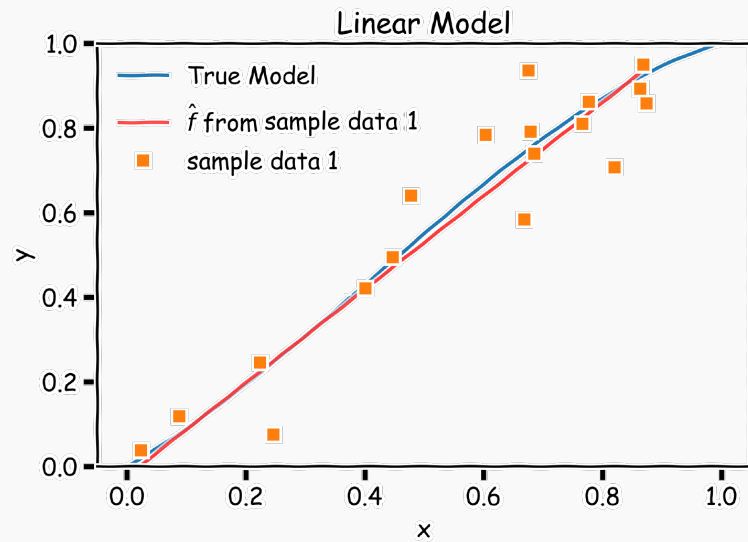Irreducible error (or aleatoric error) : we can't do anything to decrease error due to noise.

Reducible error (or epistemic error): we can decrease error due to overfitting and underfitting by improving the model.
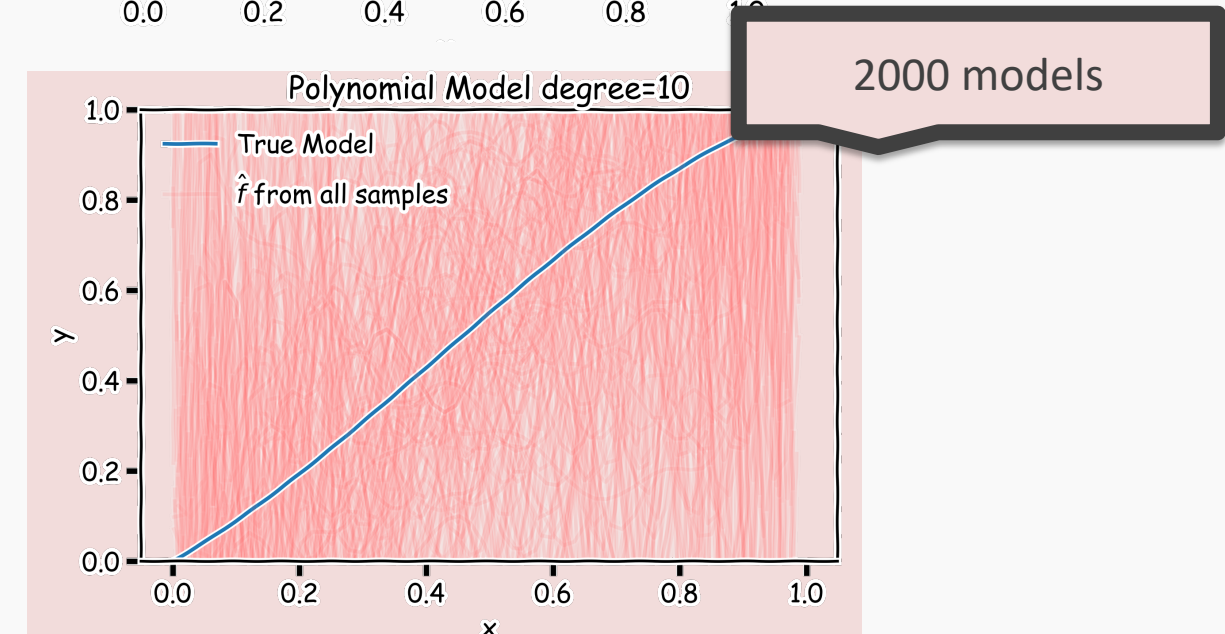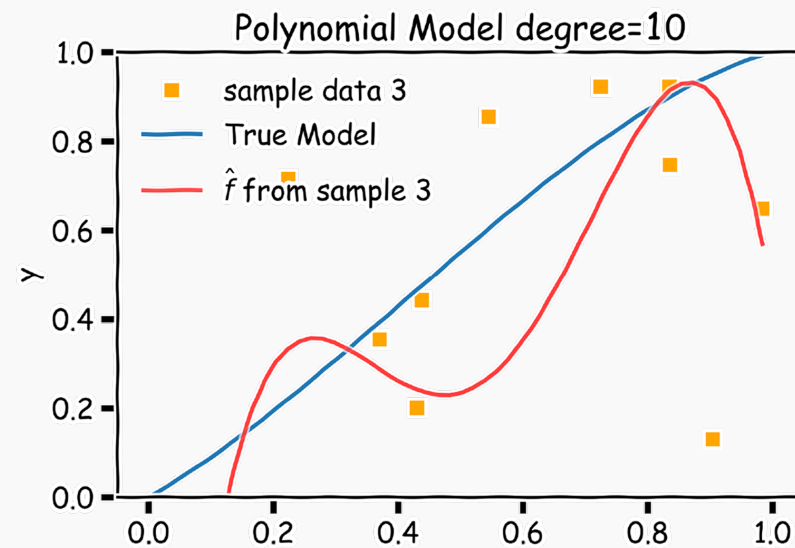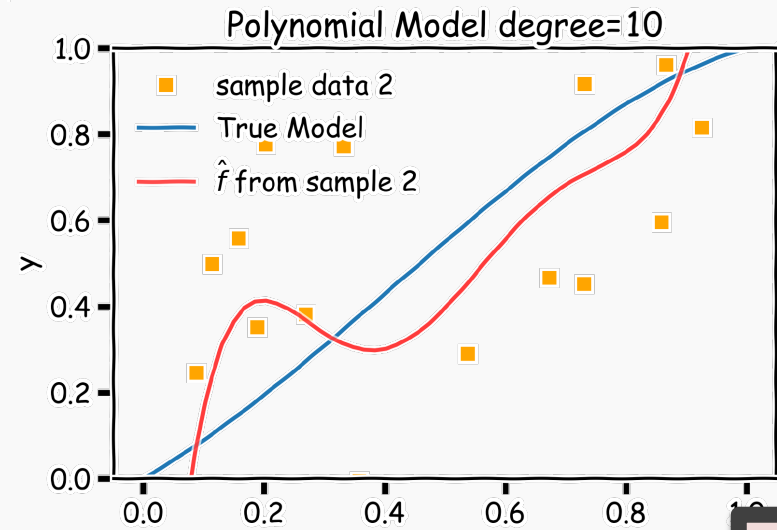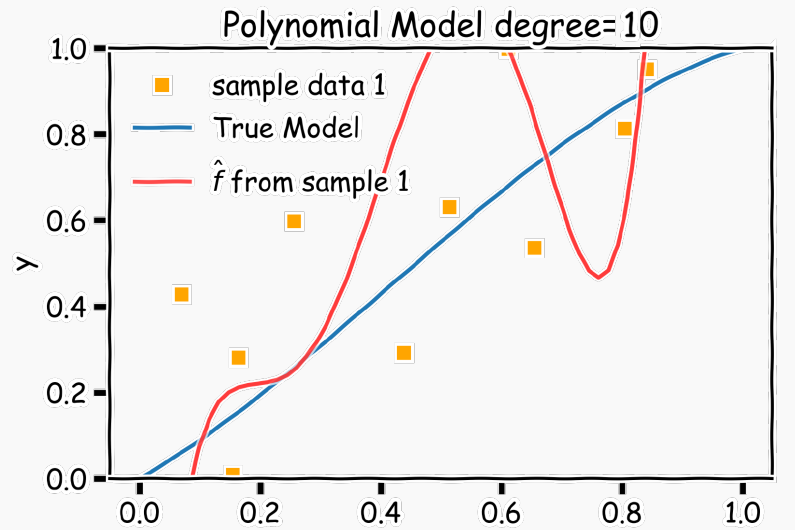
# The Bias-Variance: Bias

Reducible error comes from either underfitting or overfitting. There is a trade-off between the two sources of errors:
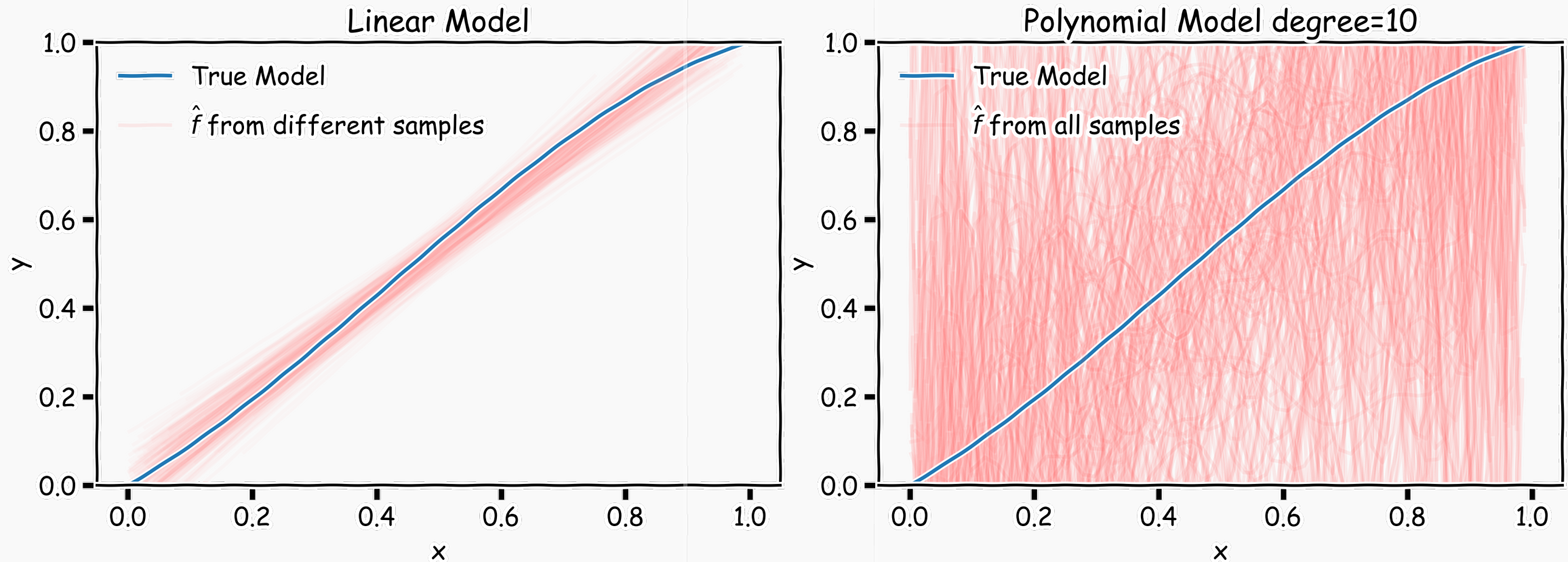
# Bias vs Variance: Variance of a SIMPLE model



2000 models

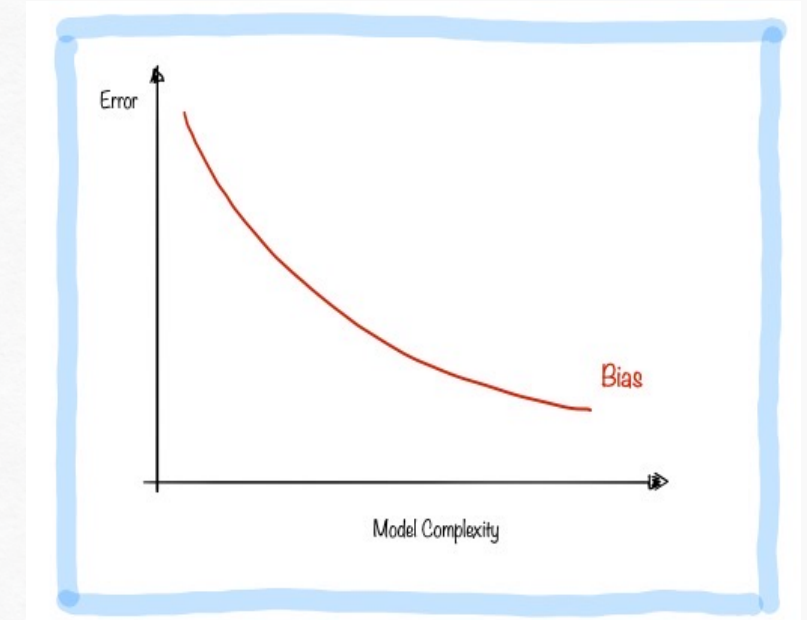# Bias vs Variance: Variance of a COMPLEX model
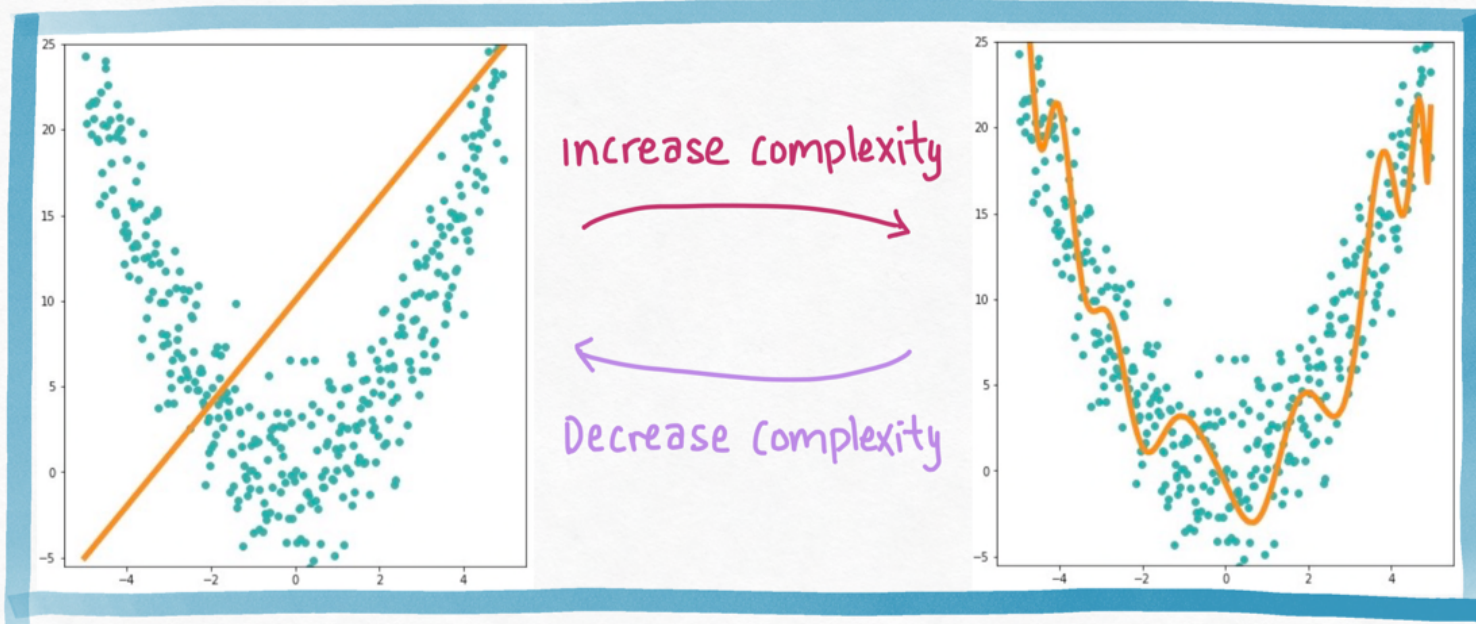
# Bias vs Variance

**Left**: 2000 best fit linear models, each fitted on a different 20-points training set.

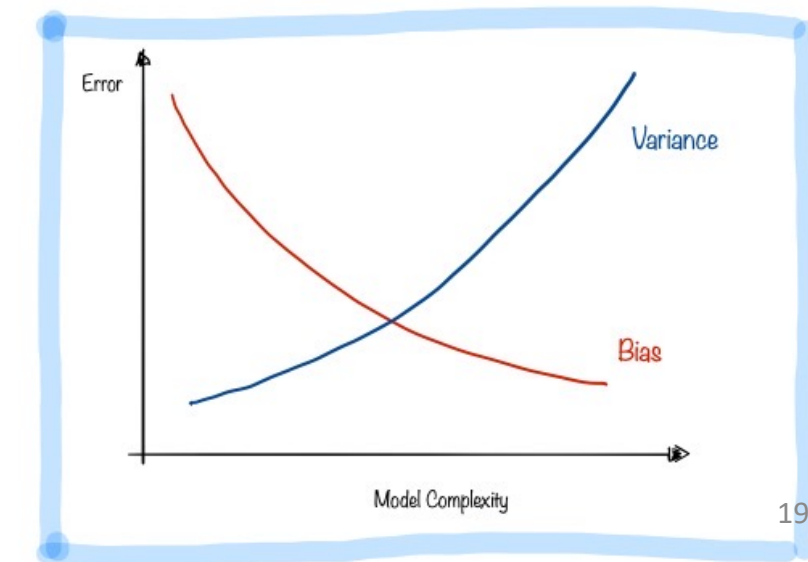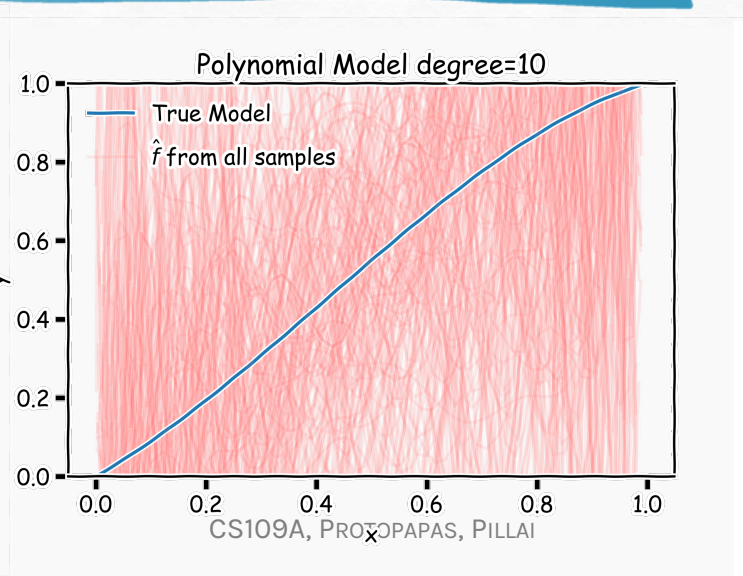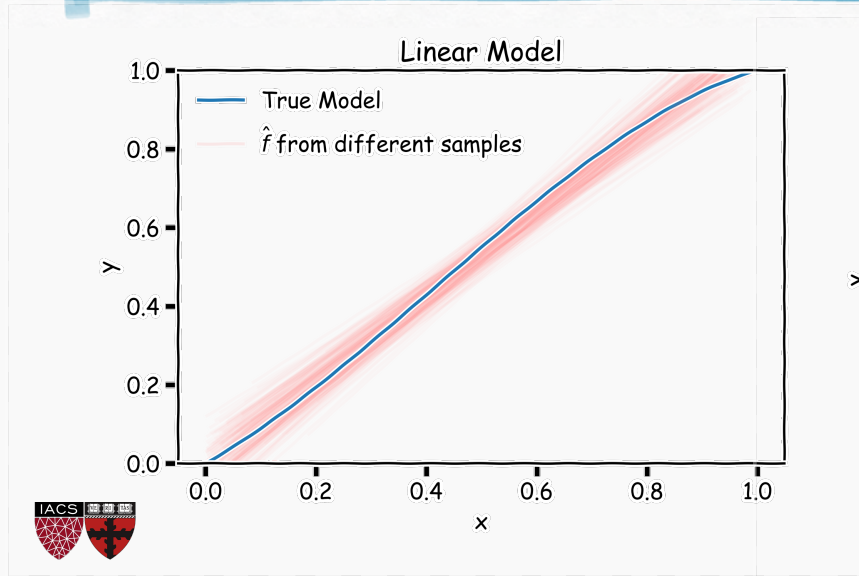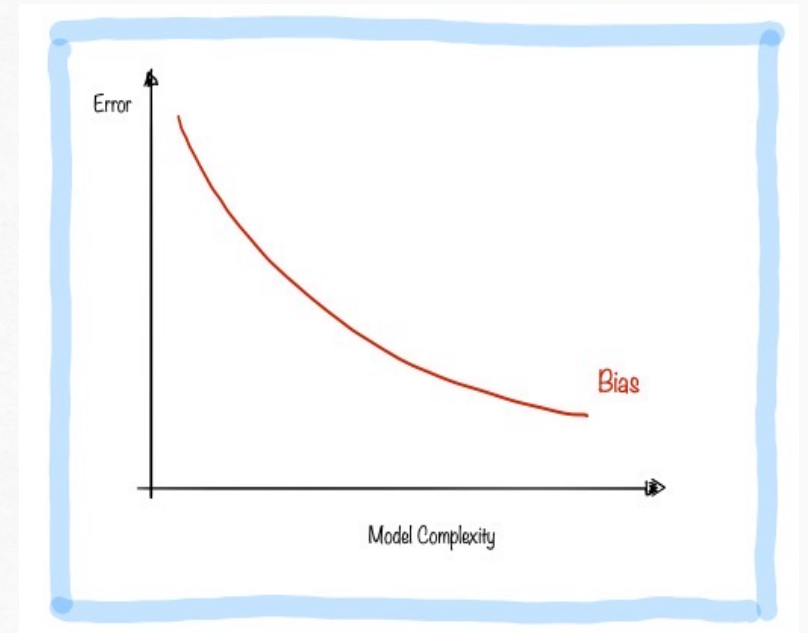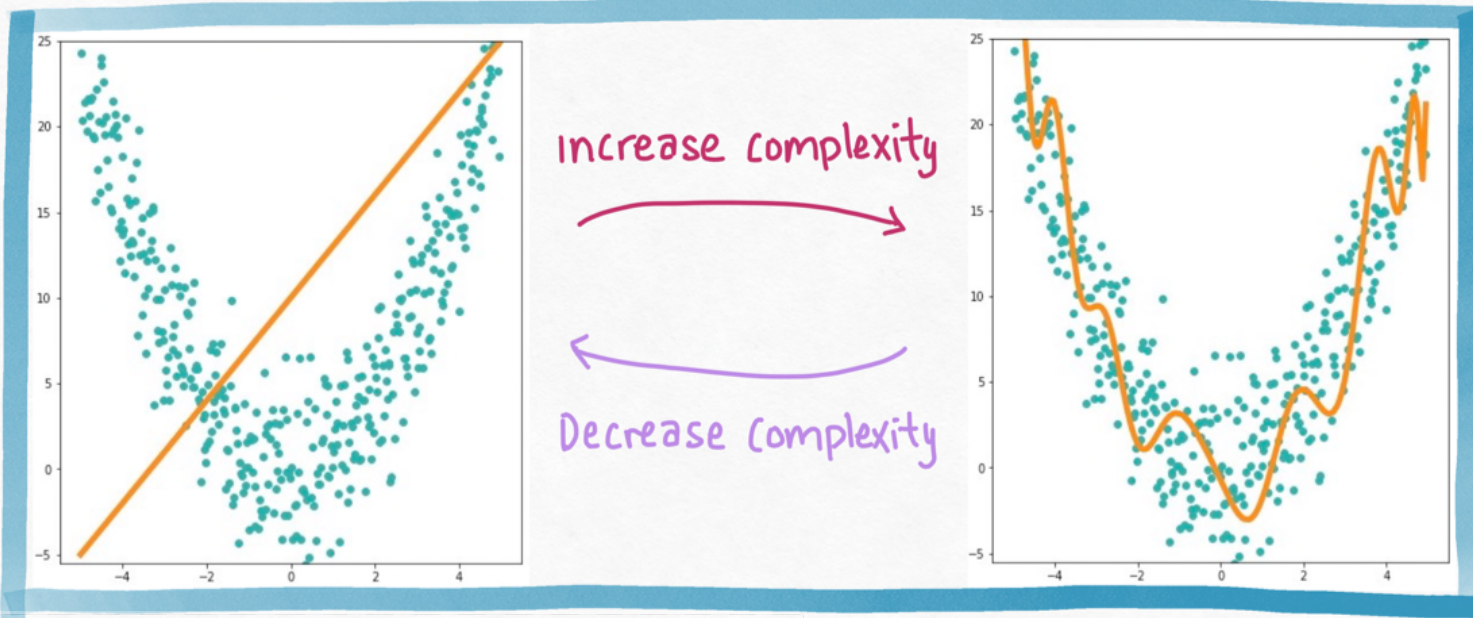**Right**: 2000 best fit models using degree 10 polynomials.

# The Bias-Variance: Bias

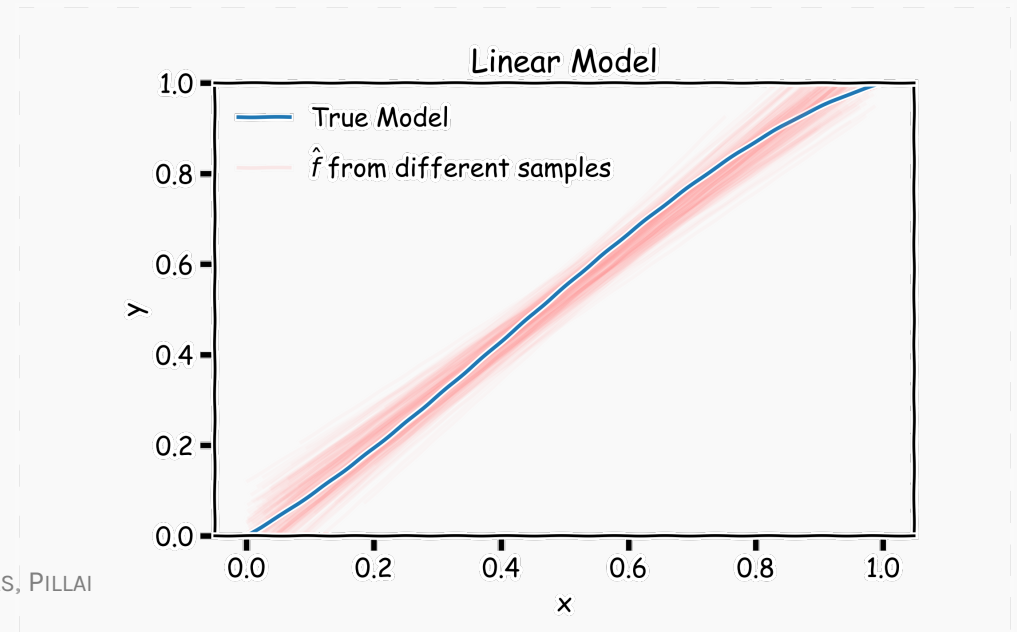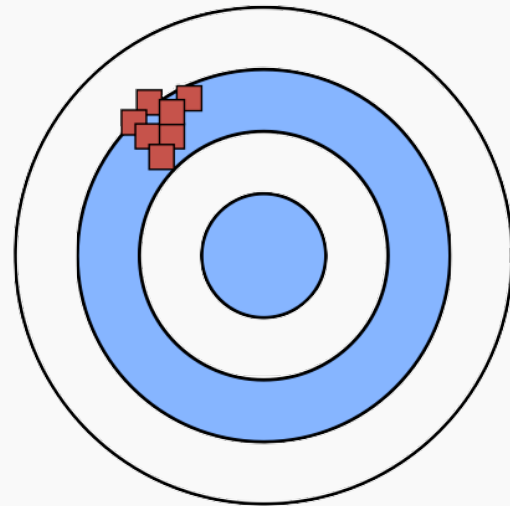Reducible error comes from either underfitting or overfitting. There is a trade-off between the two sources of errors:

# The Bias-Variance Trade Off

**Low Variance**
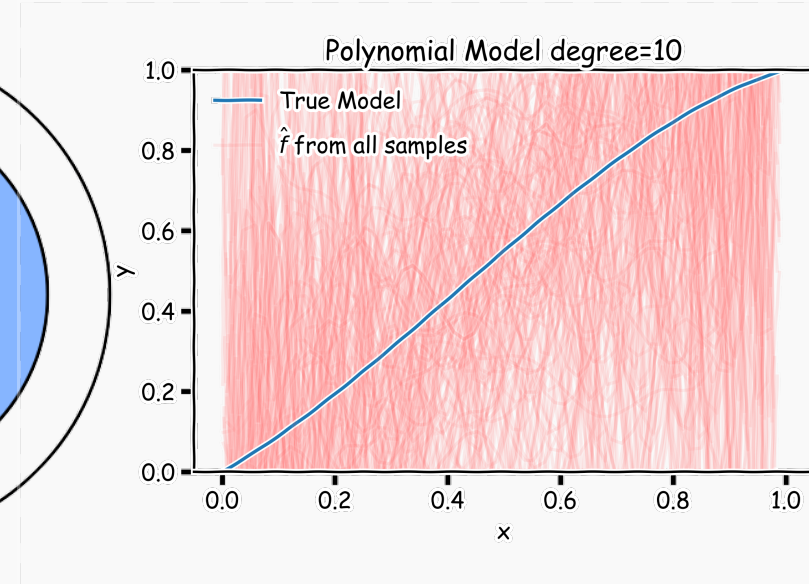(Precise)

**High Bias**
(Not Accurate)

Linear Model

True Model

$\hat{f}$ from different samples

**Low Variance**
(Precise)

**High Variance**
(Not Precise)

**Low Bias**
(Accurate)



**Polynomial Model degree=10**



**High Bias**
(Not Accurate)

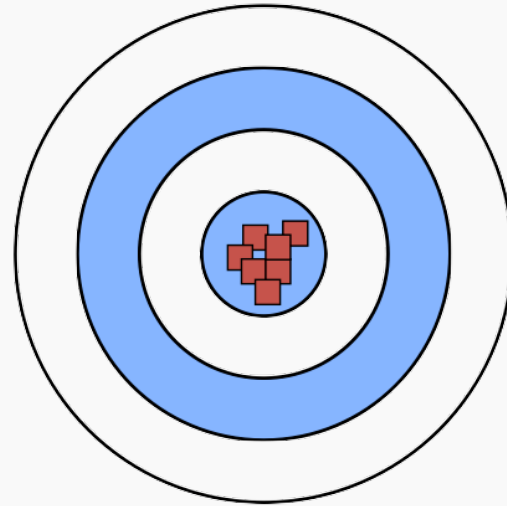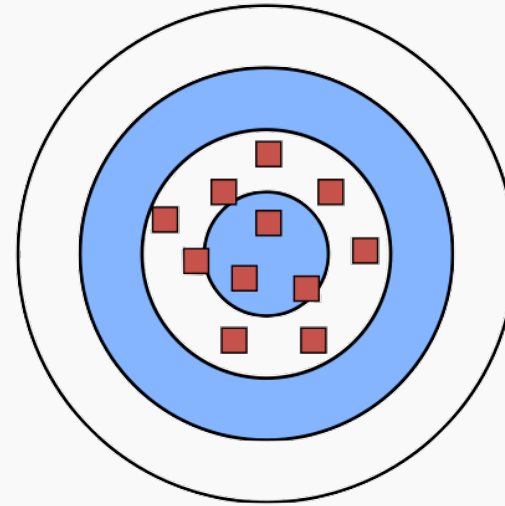**Low Variance**
(Precise)

**High Variance**
(Not Precise)

WE WANT THIS

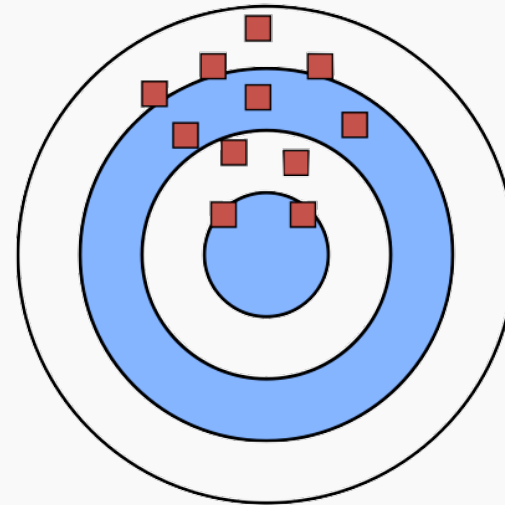**Low Bias**
(Accurate)

**High Bias**
(Not Accurate)

WE WANT TO AVOID THIS

# Overfitting

Overfitting occurs when a model corresponds too closely to the training set, and as a result, the model fails to fit additional data.

So far, we have seen that overfitting can happen when:

- Too many parameters

- Degree of the polynomial is too large

- Too many interaction terms

Soon, we will see other evidence of overfitting, which will point to a way of avoiding overfitting: Ridge and Lasso regressions.

# 👩‍🏫 Exercise: Bias Variance Tradeoff

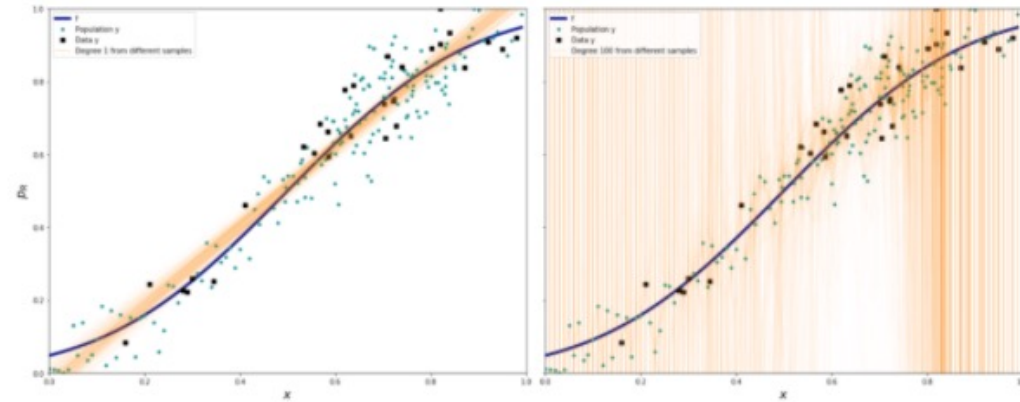The aim of this exercise is to understand **bias variance tradeoff**. For this, you will fit a polynomial regression model with different degrees on the same data and plot them as given below.



## Instructions:

- Read the file `noisypopulation.csv` as a Pandas dataframe.
- Assign the response and predictor variables appropriately as mentioned in the scaffold.
- Perform sampling on the dataset to get a subset.
- For each sampled version fo the dataset:
    - For degree of the chosen degree value:
        - Compute the polynomial features for the training
        - Fit the model on the given data
        - Select a set of random points in the data to predict the model
        - Store the predicted values as a list
- Plot the predicted values along with the random data points and true function as given above.