# Where does the error come from?

# Bias and Variance of Estimator

- Estimate the mean of a variable x
    - assume the mean of x is $\mu$
    - assume the variance of x is $\sigma^2$
- Estimator of mean $\mu$
    - Sample N points: $\{x^1, x^2, \dots, x^N\}$

$$m = \frac{1}{N}\sum_n x^n \neq \mu$$

- Estimator of variance $\sigma^2$
    - Sample N points: $\{x^1, x^2, \dots, x^N\}$

$$m = \frac{1}{N}\sum_n x^n$$
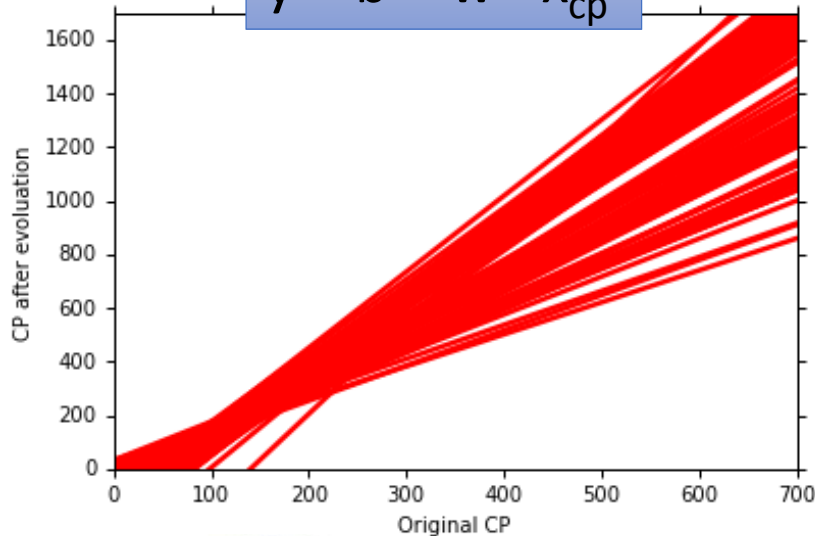
$$s = \frac{1}{N}\sum_n (x^n - m)^2$$

$$\mathrm{Var}[m] = \frac{\sigma^2}{N}$$
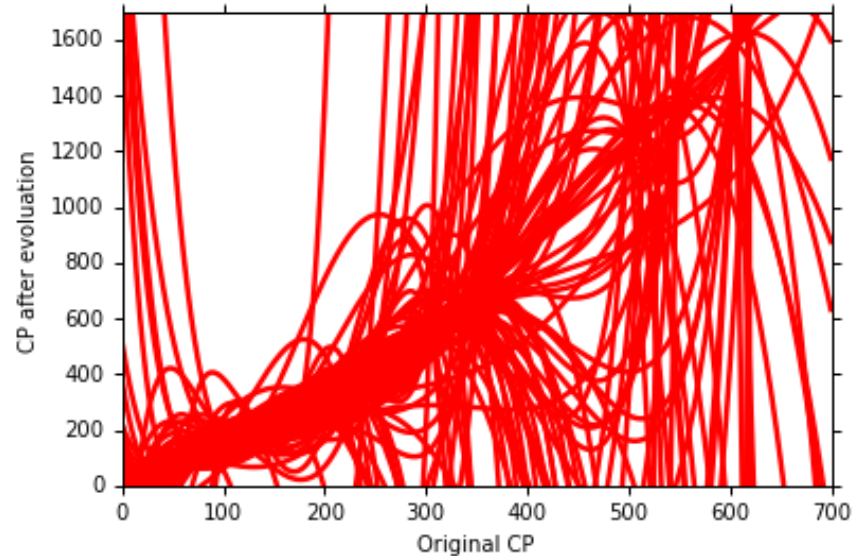
Biased estimator

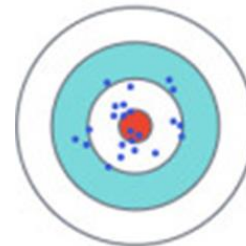$$E[s] = \frac{N-1}{N}\sigma^2 \neq \sigma^2$$

# Variance

$$y = b + w \cdot x_{cp}$$

$$y = b + w_1 \cdot x_{cp} + w_2 \cdot (x_{cp})^2 + w_3 \cdot (x_{cp})^3 + w_4 \cdot (x_{cp})^4 + w_5 \cdot (x_{cp})^5$$



Small Variance



Large Variance

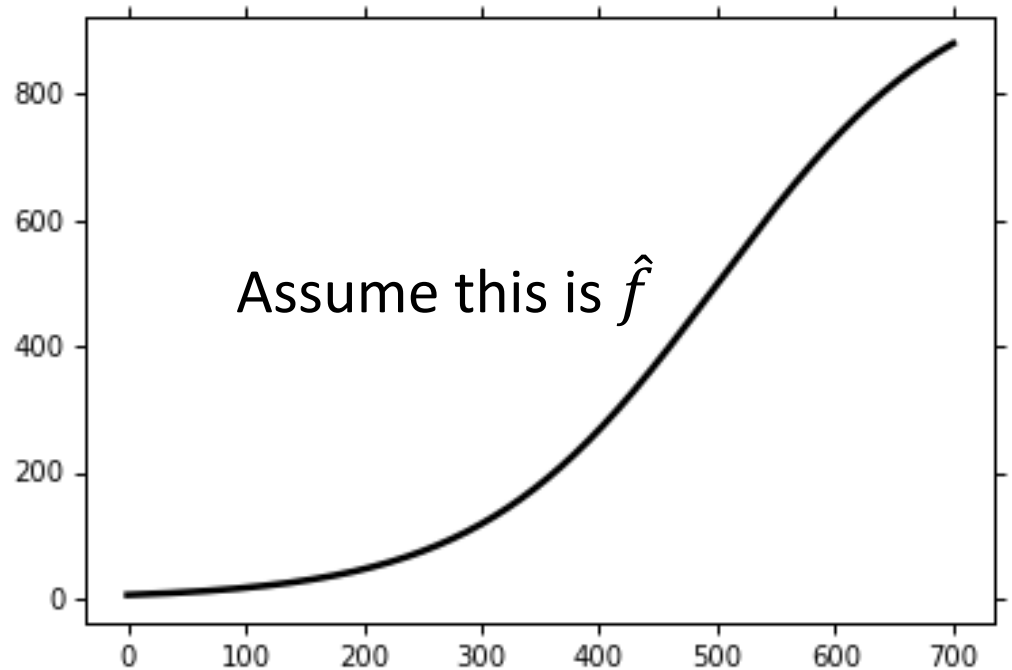Simpler model is less influenced by the sampled data

# Bias

$$E[f^*] = \bar{f}$$



Large Bias
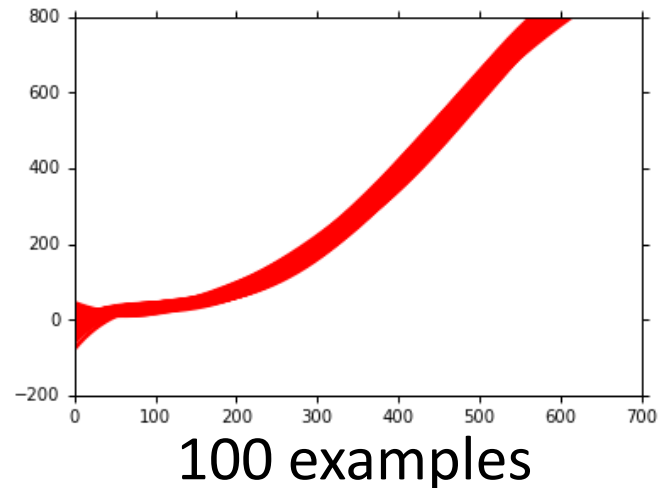
Small Bias
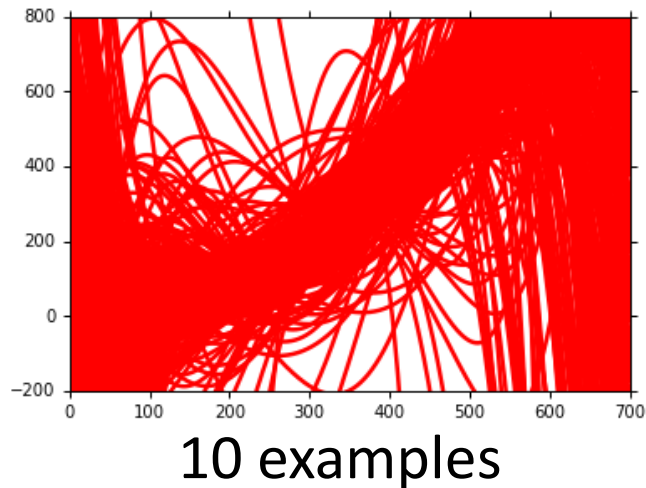
Assume this is $\hat{f}$

# What to do with large bias?

- Diagnosis:
  - If your model cannot even fit the training examples, then you have large bias   **Underfitting**
  - If you can fit the training data, but large error on testing data, then you probably have large variance   **Overfitting**

- For bias, redesign your model:
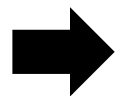  - Add more features as input
  - A more complex model
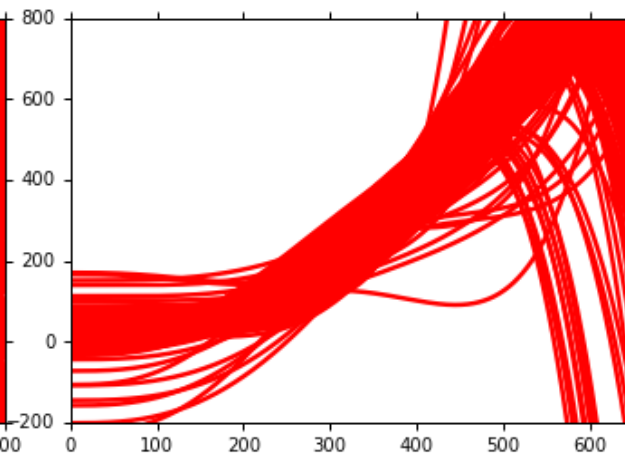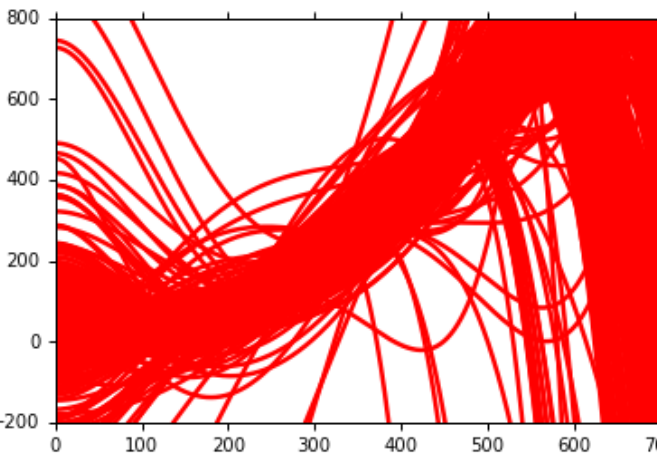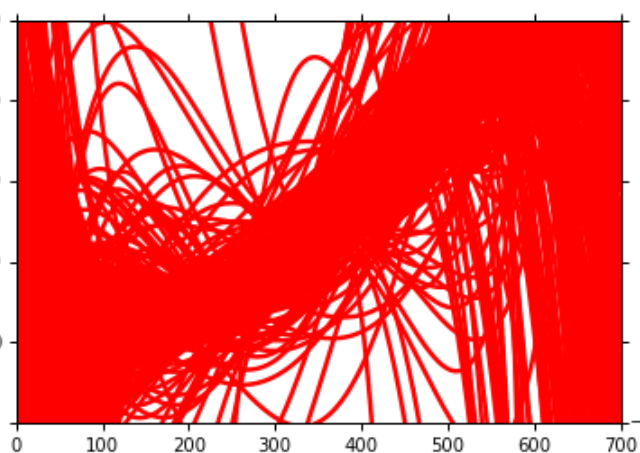
# What to do with large variance?

- More data

Very effective, but not always practical



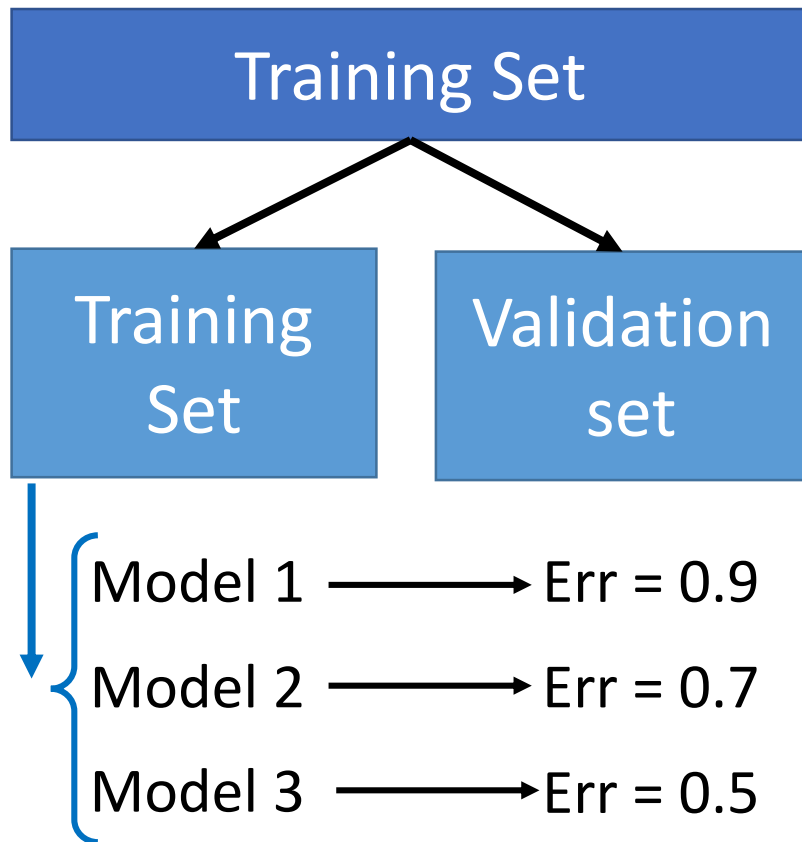10 examples



100 examples

- Regularization ➡ May increase bias

# Model Selection

- There is usually a trade-off between bias and variance.
- Select a model that balances two kinds of error to minimize total error

# Cross Validation

Training Set

private

Testing Set    Testing Set

Training Set → Validation set

Model 1 ——→ Err = 0.9

Model 2 ——→ Err = 0.7

Model 3 ——→ Err = 0.5

Using the results of public testing data to tune your model You are making public set better than private set.