Lecture 04:
*"Uncertainty, Bootstrapping"*

# An Uncertain World



Real-life test data is messy. How can we better trust this classifier?

# An Uncertain World

You have the disease.

You have the disease. I have 30% confidence that this is true.

# Uncertainty at 3 Points

Training data
$$\mathscr{D} = \left\{ \left( x_1, y_1 \right), \left( x_2, y_2 \right), \ldots, \left( x_n, y_n \right) \right\}$$

Model 1
$$\hat{y}_i = f\left( x_i, \theta \right)$$
$$L\left( y, f\left( x, \theta \right) \right)$$

Model 2
$$\hat{y}_i = g\left( x_i, \varphi \right)$$
$$L\left( y, g\left( x, \varphi \right) \right)$$

We want to know uncertainty here

Parameter estimate
$$\hat{\theta}$$

Parameter estimate
$$\hat{\varphi}$$

Prediction $\hat{y}_i = f\left( x_i, \hat{\theta} \right)$

Prediction $\hat{y}_i = g\left( x_i, \hat{\varphi} \right)$

Test data
$$\mathscr{T} = \left\{ \left( x_1, y_1 \right), \ldots \right\}$$

Test Loss $L\left( y, f\left( x_i, \hat{\theta} \right) \right)$

Compare Models

Test Loss $L\left( y, g\left( x_i, \hat{\phi} \right) \right)$

# Parameter Uncertainty

- <u>Parameter</u> =  value which summarizes data for a population; these can be expectations (*mean*) or values which describe an input-output relationship (*slope of a linear model*)

- <u>Statistic</u> = value which summarizes data from a particular sample (*i.e. sample mean*).

- <u>Estimation</u> = use a *statistic* to estimate a *parameter* of the distribution of a random variable, where
  - Estimator ($\hat{\theta}$): function used to compute <u>estimate</u>
  - Estimand ($\theta$): parameter of interest
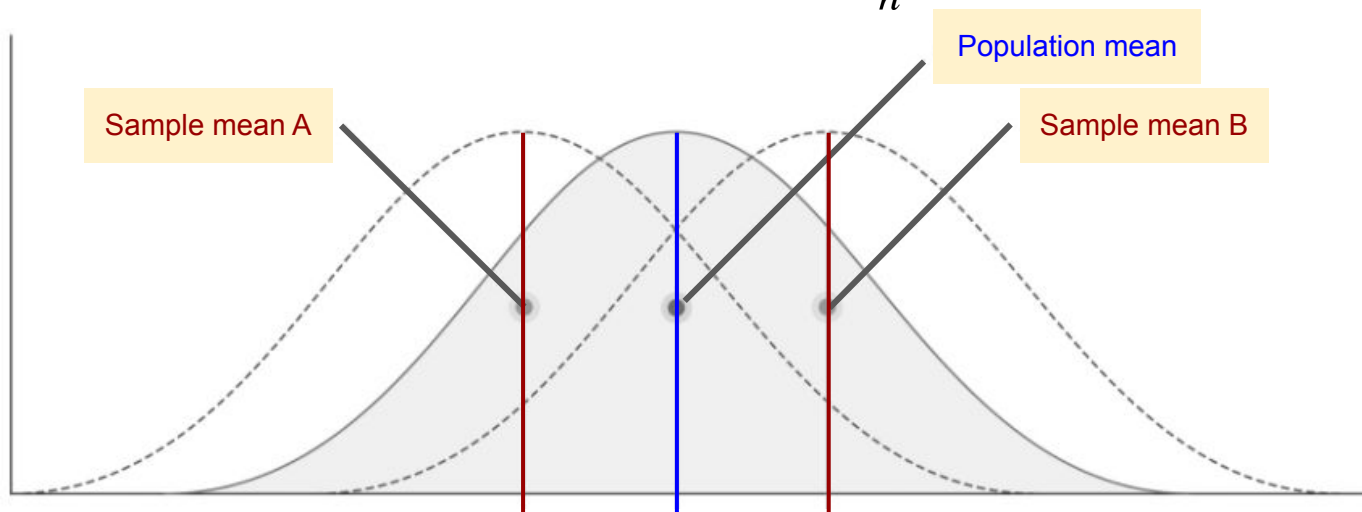
Target Population

Sample

$\mu_X$ = parameter

$\overline{x}_n$ = statistic

# Example of a parameter: mean

- Consider a model which predicts the mean… i.e. $\hat{y} = \theta$
- Given a dataset $\{x_1, x_2, \ldots, x_n\}$, the estimate for this parameter is the sample mean:

$$\widehat{\theta} = \frac{\sum_{i=1}^{n} x_i}{n}$$

Population mean

Sample mean A

Sample mean B

The distribution of an estimator is called its **sampling distribution**.

# Bias and Variance

- Bias = expected difference between estimator ($\hat{\theta}$) and parameter ($\theta$)

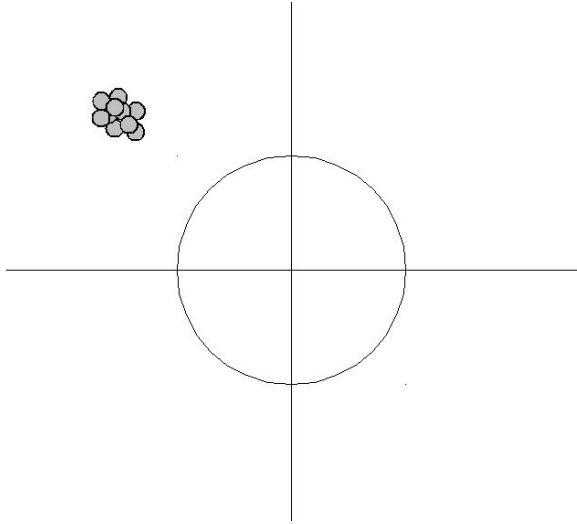  In general:  $\text{Bias}(\hat{\theta}) = E[\hat{\theta} - \theta]$

  For example:  $E[\overline{X}_n - \mu_X]$

- Variance = expected squared difference between estimator ($\hat{\theta}$) and $E[\text{estimator}]$ (mean)
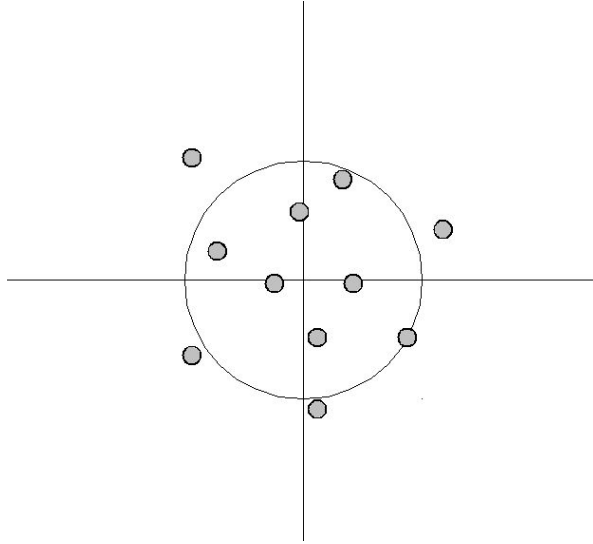
  In general:  $E\left[(\hat{\theta} - E[\hat{\theta}])^2\right]$

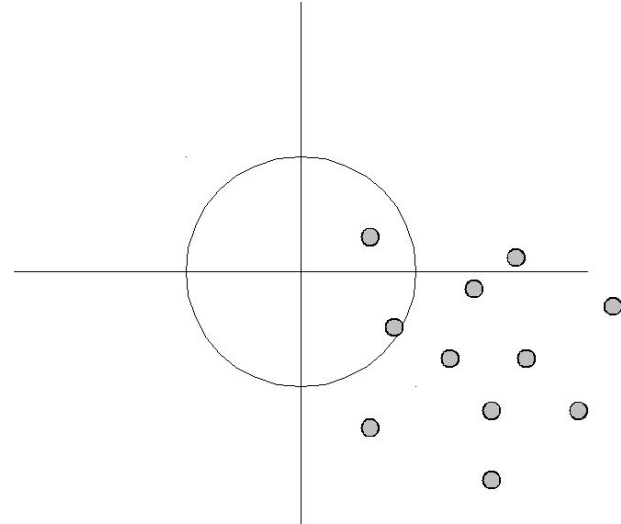  For example:  $E\left[(\overline{X}_n - E[\overline{X}_n])^2\right]$

# Bias and Variance



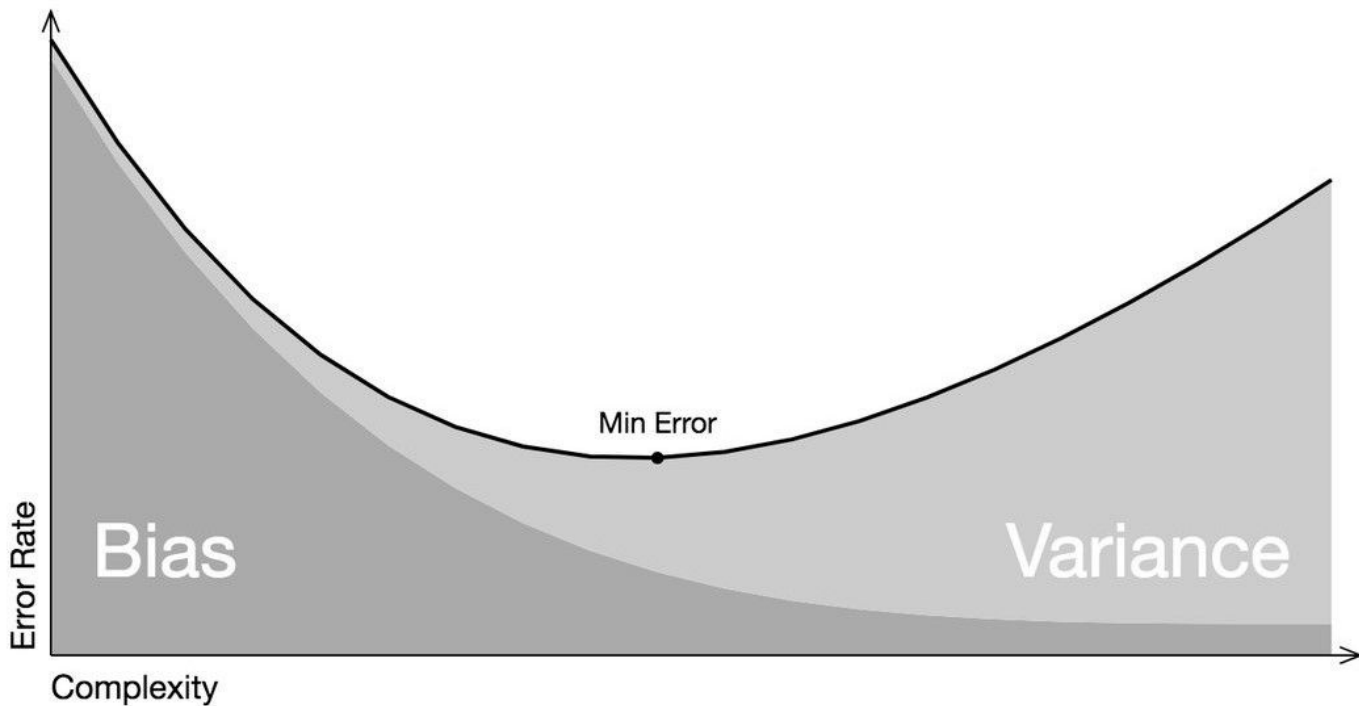High bias, low variance     Low bias, high variance     High bias, high variance

Bias-Variance Tradeoff

# Bias-Variance Tradeoff
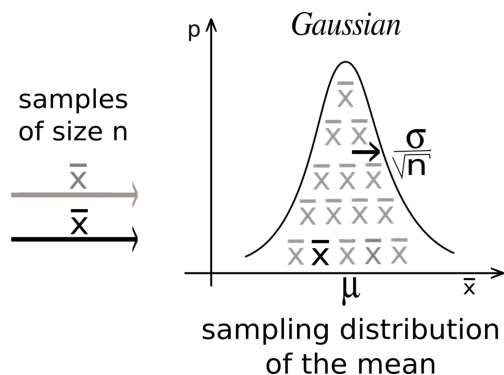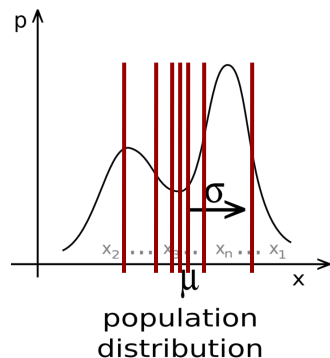
# Central Limit Theorem (CLT)

- For large n, the sampling distribution of $\overline{X}_n$ is approximately normal.
- Formally, we can write:

$$\overline{x}_n \sim N\left(\mu, \; \sigma^2 \overline{\overline{X}_n}\right), \; \text{where} \; \sigma_{\overline{X}_n} = \frac{\sigma_X}{\sqrt{n}}$$

Variance          Standard error



population distribution

samples of size n

$\overline{x}$

$\overline{x}$

*Gaussian*

$\frac{\sigma}{\sqrt{n}}$

sampling distribution of the mean

Whatever the form of the population distribution, the sampling distribution tends to a Gaussian, and its dispersion is given by the central limit theorem [1]
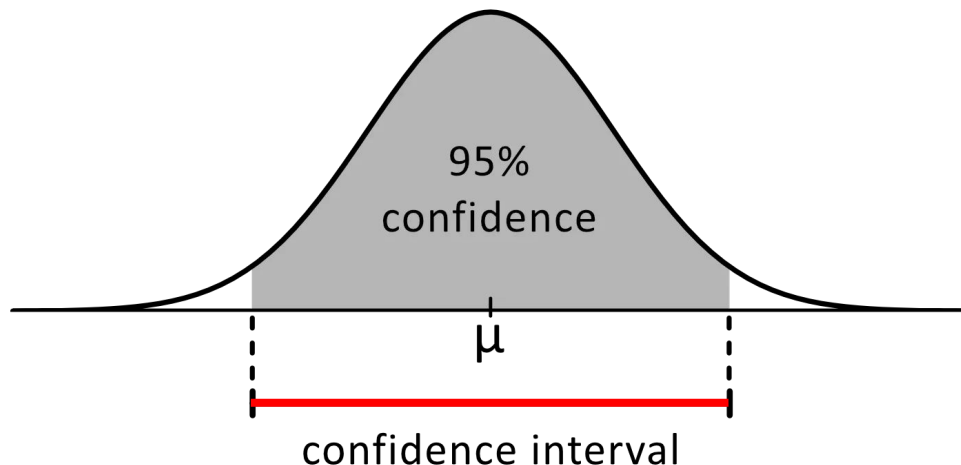
# Central Limit Theorem (CLT)

- We can use the CLT to construct **Confidence Intervals**

**Question**: What's a 95% confidence interval?

**Answer**: An interval which includes 95% of the sample means.

**Another Answer**: If we constructed this interval 100 times, it would contain the true mean in 95 of those instances.

Distribution of sample means ($\overline{x}$) around population mean ($\mu$)

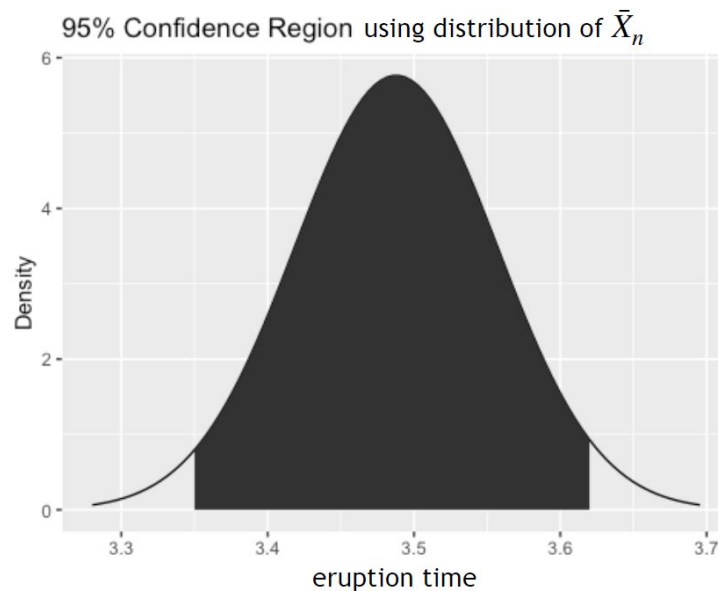95% confidence

$\mu$

confidence interval

# Central Limit Theorem (CLT)

- We can also say that 95% of the sample means are between $\mu - 1.96\sigma$ and $\mu + 1.96\sigma$
- Alternatively, 95% of the time the true mean $\mu$ will be between $\bar{x}_n - 1.96\sigma$ and $\bar{x}_n + 1.96\sigma$
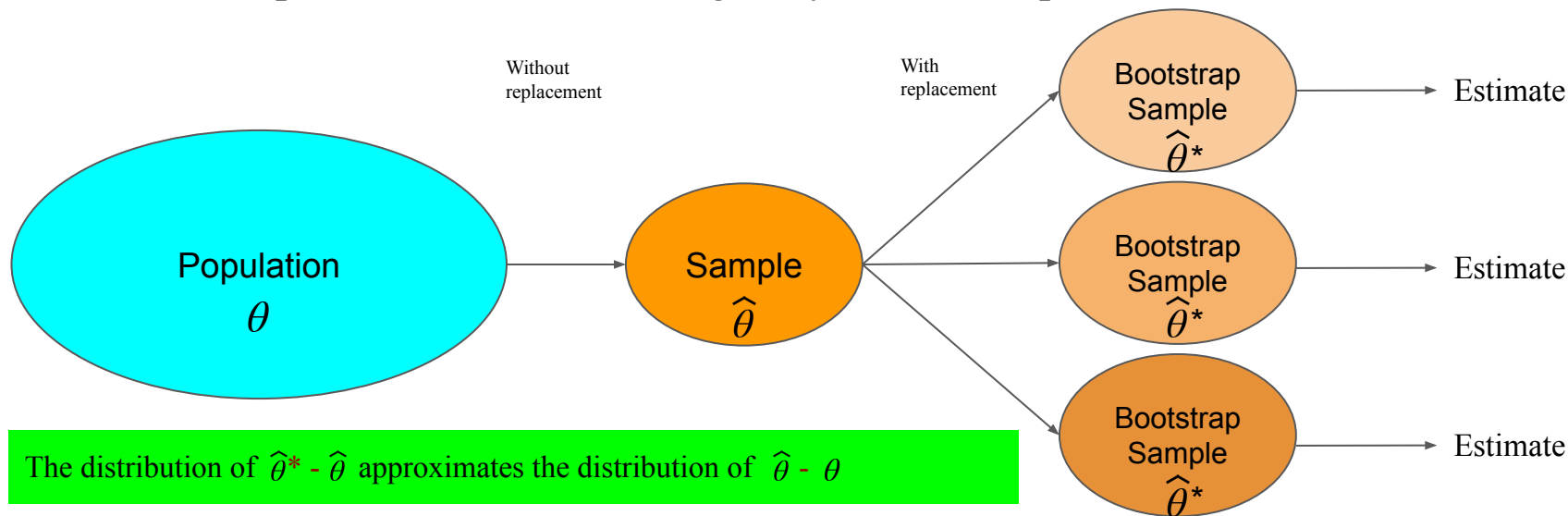
Example: Old Faithful

- Say we estimate that the mean value of eruption times is 3.4877831 (with n=272 observations)
- Is this a good estimate? How good is it?

- Mean = 3.49, Stdev = 1.14, SE = 0.07
- CI is therefore 3.49 +/- 1.96*(0.07) = 3.49 +/- 0.14



95% Confidence Region using distribution of $\bar{X}_n$

# The Bootstrap

- CLT excellent for datasets with approximately gaussian noise, and does a good job getting a distribution of parameter estimates. What if standard errors not normal?
- Bootstrap = a powerful technique to construct confidence intervals using artificially drawn samples in addition to an originally-drawn sample

Without replacement

With replacement

Population
$\theta$

Sample
$\widehat{\theta}$

Bootstrap Sample
$\widehat{\theta}*$

Estimate

Bootstrap Sample
$\widehat{\theta}*$

Estimate

Bootstrap Sample
$\widehat{\theta}*$

Estimate

The distribution of $\widehat{\theta}*$ - $\widehat{\theta}$ approximates the distribution of $\widehat{\theta}$ - $\theta$

# The Bootstrap

```
Your Sample S has N observations
For b in 1:numBootstrap:
    resample N from S with replacement -> S*
    Fit model to S* -> θ̂*(bootstrap statistics)
    Record your bootstrap statistics

Return the distribution of bootstrap statistics
```
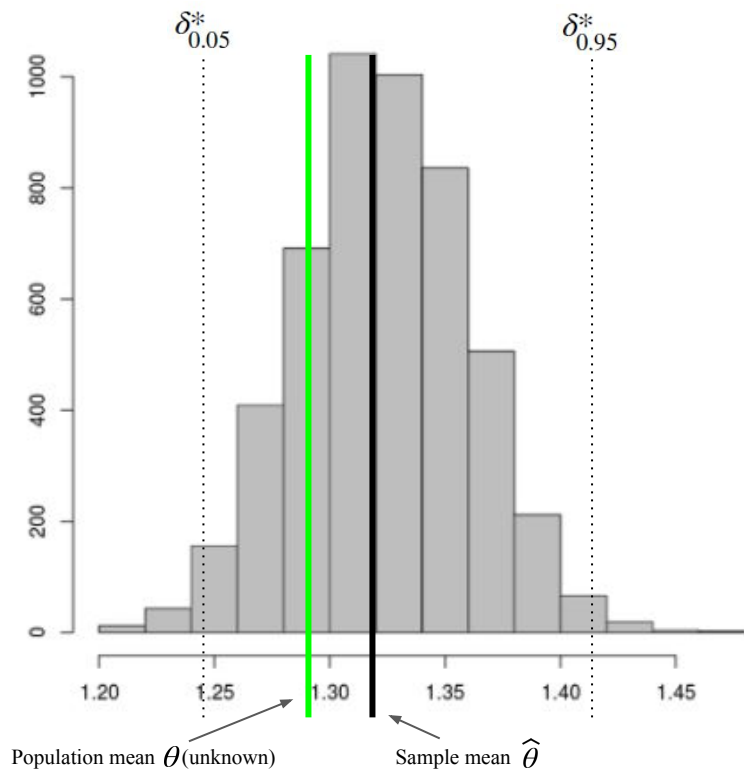
# The Bootstrap

- Now that we have a distribution of bootstrap statistics, we can construct a CI



Population mean $\theta$ (unknown)    Sample mean $\widehat{\theta}$

- For example, a 90% confidence interval centred at the sample mean would be

$$CI = \left[ \hat{\theta} - \delta^*_{0.95}, \hat{\theta}^* - \delta^*_{0.05} \right]$$

where $\quad \delta^* = \hat{\theta}^* - \hat{\theta}$

and where $\delta^*_{0.95}$ is the 95% percentile of the bootstrap distribution

# Prediction Uncertainty

Training data
$$\mathcal{D} = \left\{ \left( x_1, y_1 \right), \left( x_2, y_2 \right), \ldots, \left( x_n, y_n \right) \right\}$$

Model
$$\widehat{y}_i = f\left( x_i, \theta \right)$$

$$L\left( y, f\left( x, \theta \right) \right)$$
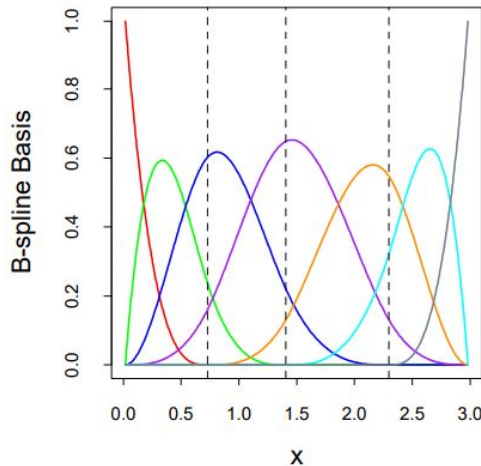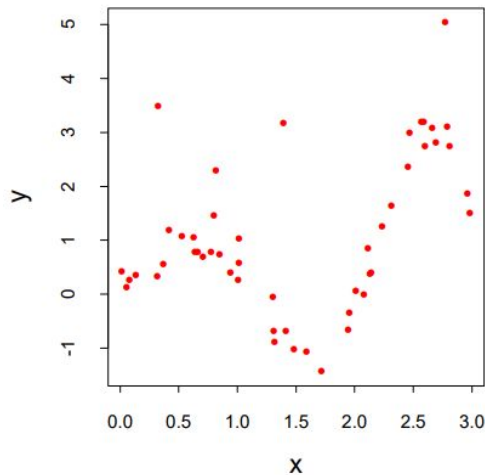
Parameter estimate
$$\widehat{\theta}$$

Prediction
$$\widehat{y}_i = f\left( x_i, \widehat{\theta} \right)$$

How does uncertainty in our parameter estimate influence the uncertainty of our prediction?

How much would the prediction change if we had used a different set of training data?
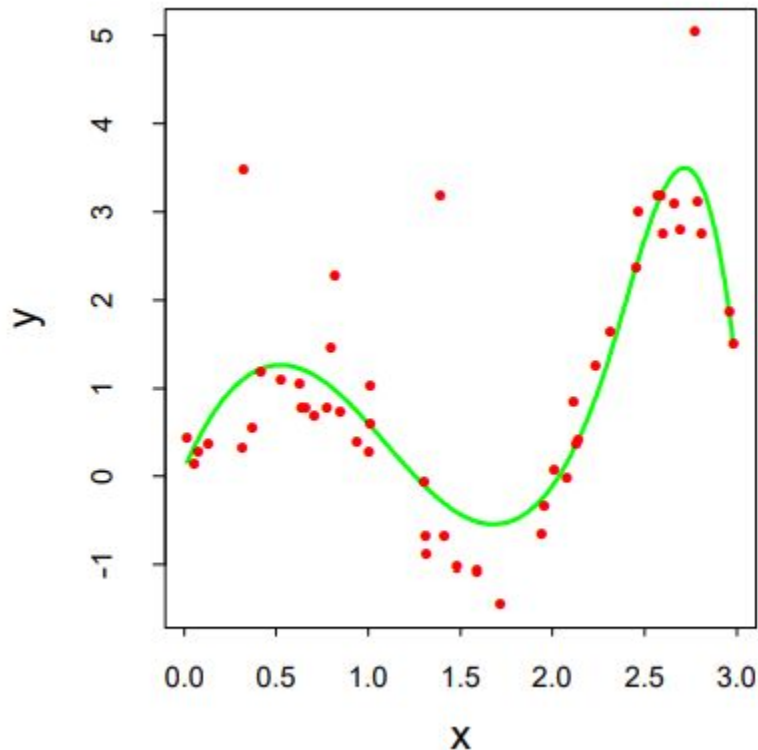
# Prediction Uncertainty (Bootstrap)



- Example: say we want to fit a cubic spline to this data. We can use a linear expansion of B-spline basis functions $h_i(x)$.

- We store the B coefficients of these basis functions into a vector $\theta$, and fit $\hat{y} = f(x) = X\theta$.

$$\begin{pmatrix} y_1 \\ y_2 \\ \cdots \\ y_n \end{pmatrix}_{n \times 1} = \begin{pmatrix} h_1(x_1) & h_2(x_1) & \cdots & h_p(x_1) \\ h_1(x_2) & h_2(x_2) & \cdots & h_p(x_2) \\ & & & \\ h_1(x_n) & h_2(x_n) & \cdots & h_p(x_n) \end{pmatrix}_{n \times p} \begin{pmatrix} \beta_1 \\ \cdots \\ \beta_p \end{pmatrix}_{p \times 1}$$

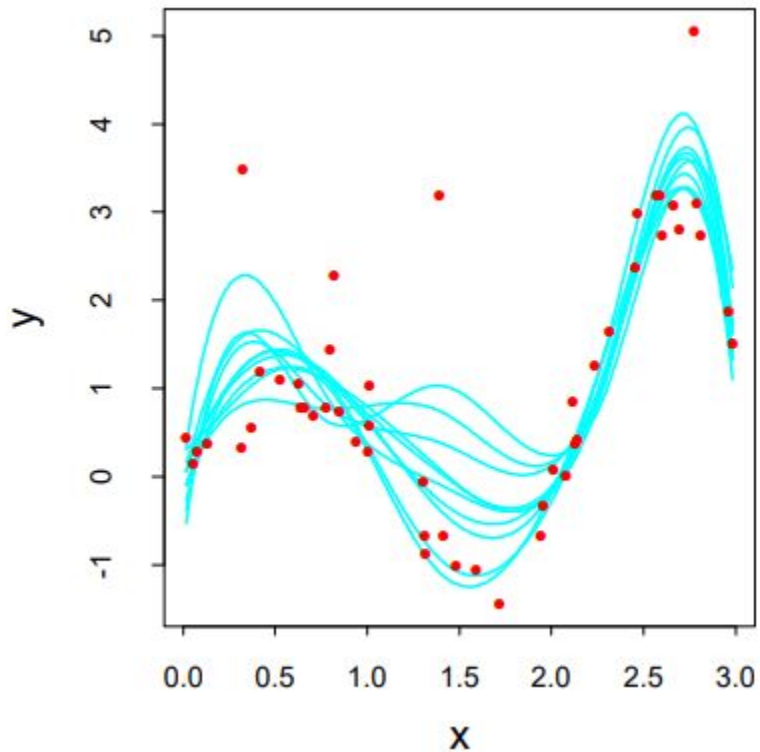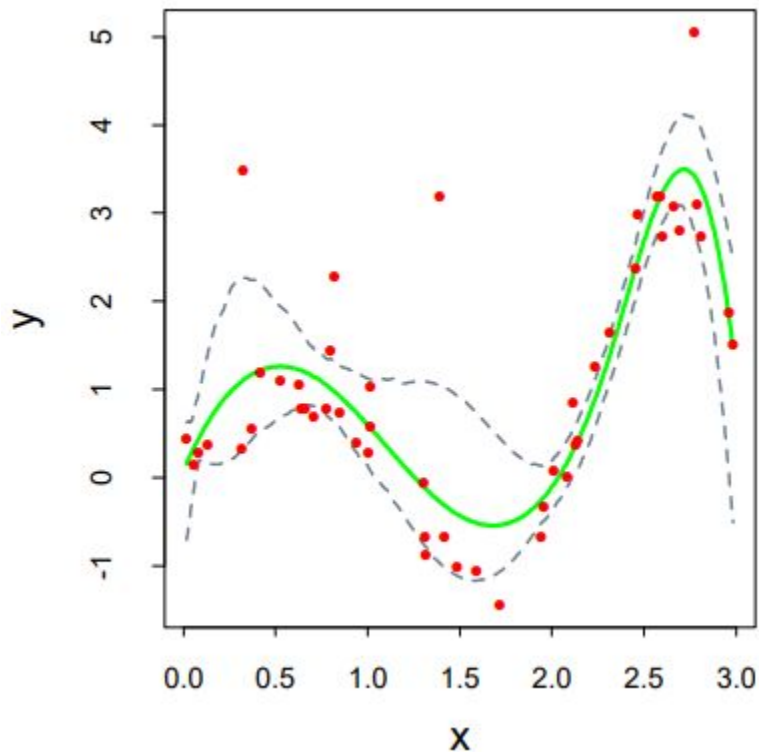[†] See ESLII, page 262

# Prediction Uncertainty (Bootstrap)



- Here is our fit $\hat{y}=\widehat{f}(x)=X\widehat{\theta}$

- Is it any good? Yes, but how confident can we be of this?

- **Let's use bootstrap**:
  - From our original sample, generate a new sample (with replacement)
  - For this new sample, get a new parameter estimate $\widehat{\theta}_b^*$
  - Do this as many times as you can

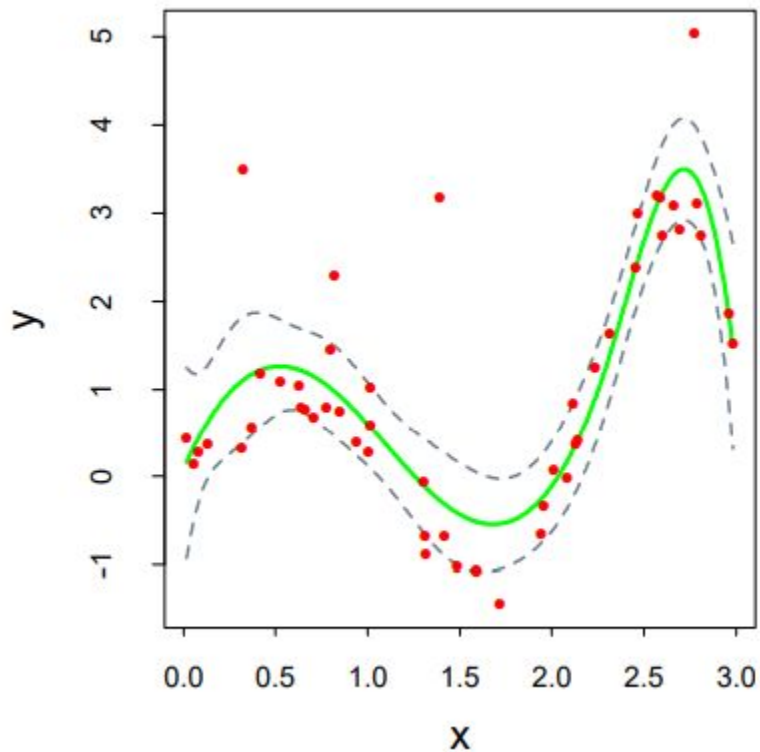# Prediction Uncertainty (Bootstrap)



- We can plot each new prediction $\hat{y}_b^* = \hat{f}_b^*(x) = X\hat{\theta}_b^*$
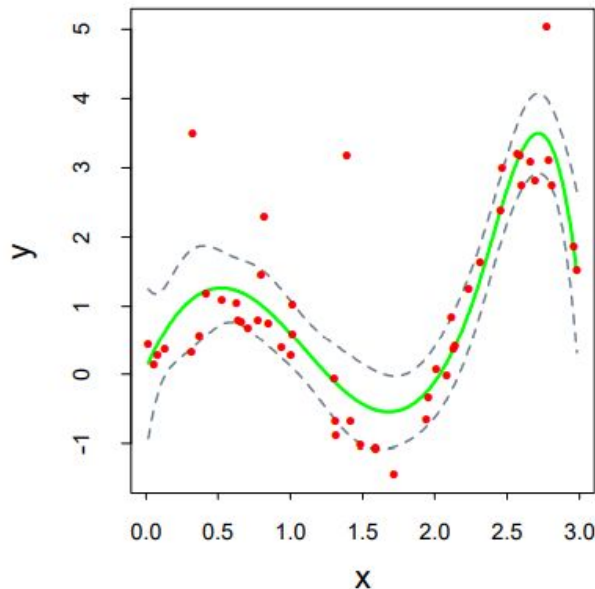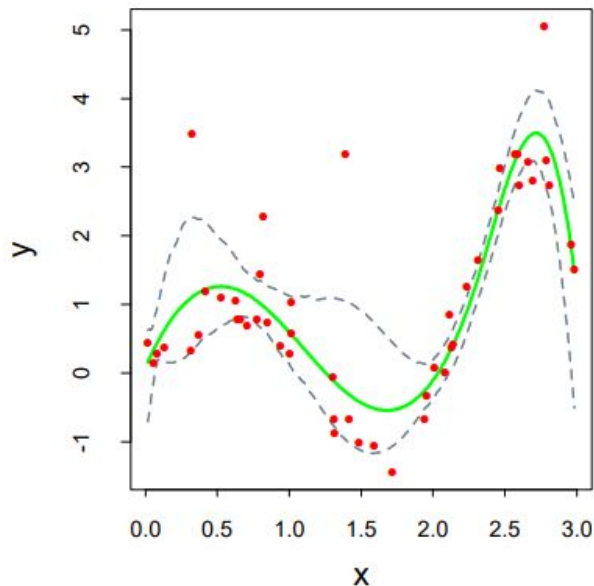
# Prediction Uncertainty (Bootstrap)



- And since we now have a distribution of samples for each x, we can compute a 95% Confidence Interval (CI)
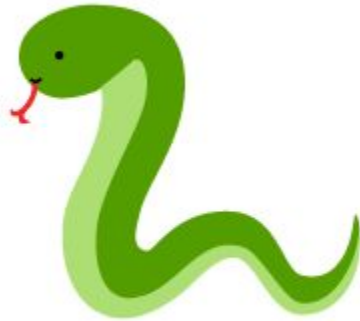
# Prediction Uncertainty (Bootstrap)



- Note that we could have also used CLT to get the CIs

# Prediction Uncertainty (Bootstrap)



- <u>Warning</u>: our Confidence Interval (via bootstrap [left] or CLT [right]) is for the true value of $f(x)$, not for new observations $(x_{new}, y_{new})$
- Why? A CI for new data would need to also consider random variability ($\sigma^2$) between $f(x_n)$ and $y_n$.

Let'sss try it in Python...

# Summary

- Parameter Uncertainty
  - Parameters, Statistics, Estimation
  - Example using Population/Sample Mean
  - Bias and variance
  - The Central Limit Theorem (CLT)
  - Constructing a Confidence Interval (CI)
  - Bootstrap
- Prediction Uncertainty
  - B-spline example (Bootstrap)
- Coding examples of CLT, Bootstrap