# Agenda

1. Interpretation with transformed variables

2. Visualizing model predictions

3. Examining model fit

4. Overfitting and underfitting

5. Visualizing predictions in R

# Interpreting coefficients

Predicting:

**Income** ⟵ **Age**

Using:

$$\text{Income}_i \sim \text{Norm}(\mu_i, \sigma)$$
$$\mu_i = \alpha + \beta \text{Age}_i$$

| | Post. Mean |
|---|---|
| $\alpha$ | 32813.3 |
| $\beta$ | 185.7 |

$$E(\text{Inc.}|a_2) - E(\text{Inc.}|a_1) = (\alpha + \beta a_2) - (\alpha + \beta a_1)$$
$$= \beta(a_2 - a_1)$$
$$= 185.7(a_2 - a_1)$$

# Interpreting coefficients

Predicting:

**Income**

Using:

**Age**

$$\text{Income}_i \sim \text{Norm}(\mu_i, \sigma)$$
$$\mu_i = \alpha + \beta \text{Age}_i$$

| | Post. Mean |
| --- | --- |
| $\alpha$ | 32813.3 |
| $\beta$ | 185.7 |

**Units of age** | Years

**Units of income** | Dollars

**Interpreting $\beta$** | For each year of age, the model predicts about $186 more income per year.

# Standardized variables

**Predicting:**

## St(Income) ← **Using:** ## St(Age)

$$St(Income_i) \sim Norm(\mu_i, \sigma)$$
$$\mu_i = \alpha + \beta St(Age_i)$$

| | Post. Mean |
|---|---|
| $\alpha$ | 0 |
| $\beta$ | 0.065 |

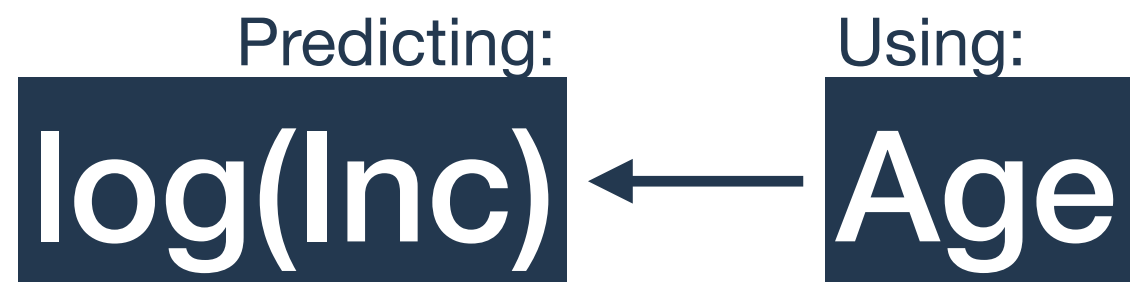**Units of age** | Standard deviations of age

**Units of income** | Standard deviations of income

**Interpreting $\beta$** | For each standard deviation of age, the model predicts an increase of about 0.065 standard deviations in income.

# Log of outcome variable

Predicting:
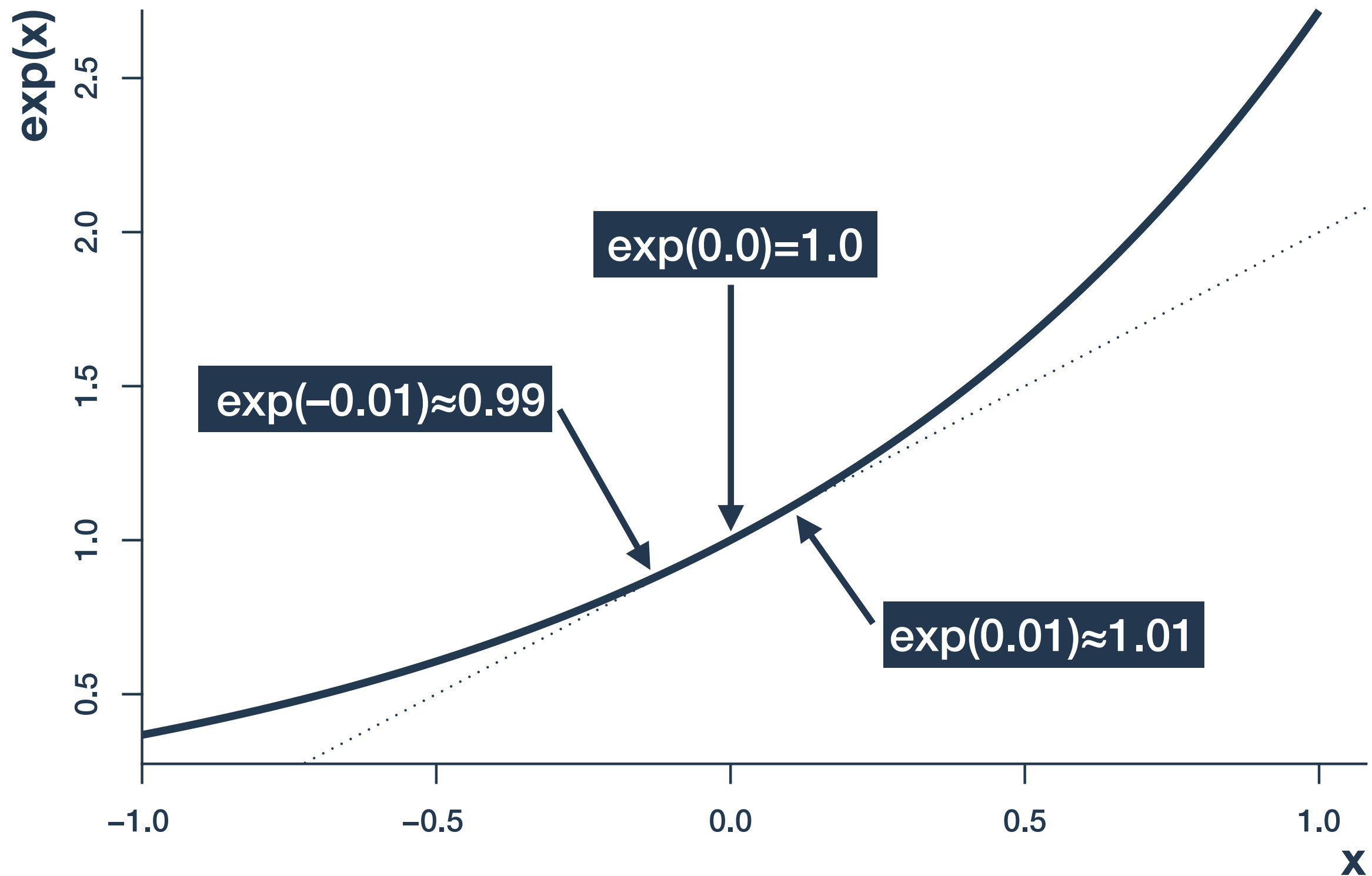
Using:

log(Inc) $\longleftarrow$ Age

$$\log(\text{Income}_i) \sim \text{Norm}(\mu_i, \sigma)$$
$$\mu_i = \alpha + \beta \text{Age}_i$$

|   | Post. Mean | exp(Mean) |
|---|---|---|
| $\alpha$ | 9.648 | 15489.124 |
| $\beta$ | 0.009 | 1.009 |

$$\text{Inc}_2 / \text{Inc}_1 = \exp(log(\text{Inc}_2) - \log(\text{Inc}_1))$$
$$= \exp((\alpha + \beta a_2) - (\alpha + \beta a_1))$$
$$= \exp(\beta(a_2 - a_1))$$

# Log of outcome variable

# Log of outcome variable

Predicting:    Using:

**log(Inc)** ← **Age**

$$\log(\text{Income}_i) \sim \text{Norm}(\mu_i, \sigma)$$
$$\mu_i = \alpha + \beta\text{Age}_i$$

|   | Post. Mean | exp(Mean) |
|---|---|---|
| $\alpha$ | 9.648 | 15489.124 |
| $\beta$ | 0.009 | 1.009 |

| | |
|---|---|
| **Units of age** | Years |
| **Units of income** | Log dollars |
| **Interpreting $\beta$** | For each year of age, the model predicts a 0.9% increase in income. |

# Log of predictor variable

Predicting:     Using:

**Income** ⟵ **log(Age)**

| | Income$_i \sim$ Norm$(\mu_i, \sigma)$ |
| --- | --- |
| | $\mu_i = \alpha + \beta\log(\text{Age}_i)$ |

| | Post. Mean |
| --- | --- |
| $\alpha$ | −17586 |
| $\beta$ | 15675 |

$$\text{Inc}_2 - \text{Inc}_1 = (\alpha + \beta\log(a_2)) - (\alpha + \beta\log(a_1)))$$
$$= \beta(\log(a_2) - \log(a_1)))$$
$$= \beta(\log(a_2/a_1)))$$

# Log of predictor variable

Predicting:     Using:

**Income** ⟵ **log(Age)**

$$\text{Income}_i \sim \text{Norm}(\mu_i, \sigma)$$
$$\mu_i = \alpha + \beta \log(\text{Age}_i)$$

|  | Post. Mean |
|---|---|
| $\alpha$ | −17586 |
| $\beta$ | 15675 |

| | |
|---|---|
| **Units of age** | Log years |
| **Units of income** | Dollars |
| **Interpreting $\beta$** | For each 10% increase in age, the model predicts an increase of $\beta \times \log(1.1) = 15{,}675 \times 0.095 = 1{,}494.07$ dollars in income. |

# Visualizing predictions

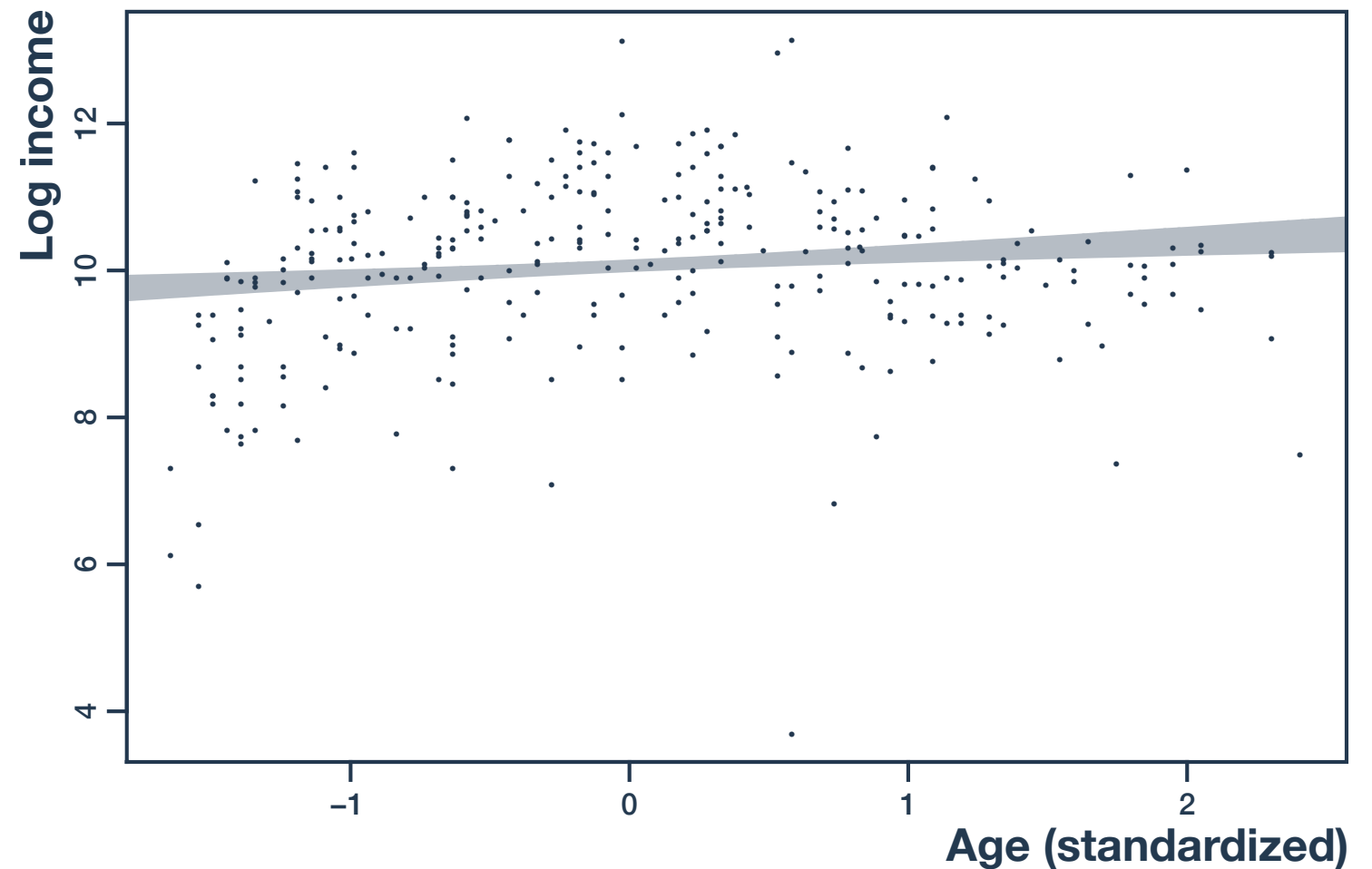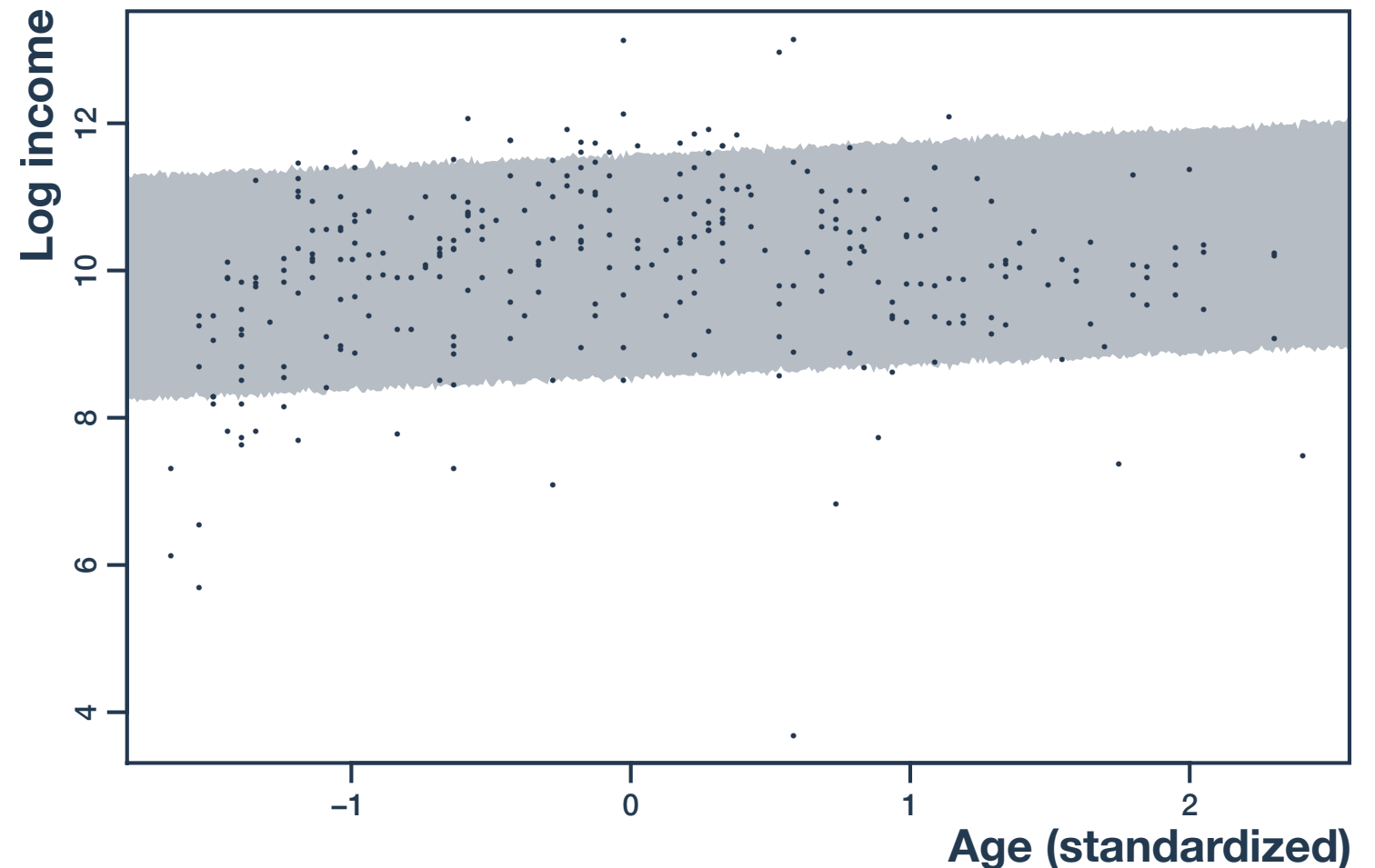$$\log(\text{Inc}_i) \sim \text{Norm}(\mu_i, \sigma)$$

$$\mu_i = \alpha + \beta \text{St}(\text{Age}_i)$$

$$\alpha \sim \text{Norm}(10, 2)$$

$$\beta \sim \text{Norm}(0, 3)$$

$$\sigma \sim \text{Unif}(0, 5)$$

# Visualizing predictions

$$\log(\text{Inc}_i) \sim \text{Norm}(\mu_i, \sigma)$$

$$\mu_i = a + \beta \text{St}(\text{Age}_i)$$

$$a \sim \text{Norm}(10, 2)$$

$$\beta \sim \text{Norm}(0, 3)$$

$$\sigma \sim \text{Unif}(0, 5)$$



| | Post. Mean | exp(Mean) |
|---|---|---|
| $a$ | 10.06 | 0.07 |
| $\beta$ | 0.17 | 0.07 |
| $\sigma$ | 1.18 | 0.05 |

**Posterior distribution of mean:**

$$\text{Pr}(\mu | \text{Age} = a)$$

1. Take a sample of size N from posterior $\text{Pr}(a, \beta, \sigma | D)$.
2. For each value of Age $a$, calculate N values of $\mu = a + \beta a$.
3. Calculate quantiles (say, 10% and 90%) for posterior of $\mu$ at each value of $a$.

# Visualizing predictions

$$\log(\text{Inc}_i) \sim \text{Norm}(\mu_i, \sigma)$$

$$\mu_i = \alpha + \beta \text{St}(\text{Age}_i)$$

$$\alpha \sim \text{Norm}(10, 2)$$

$$\beta \sim \text{Norm}(0, 3)$$

$$\sigma \sim \text{Unif}(0, 5)$$



| | Post. Mean | exp(Mean) |
|---|---|---|
| $\alpha$ | 10.06 | 0.07 |
| $\beta$ | 0.17 | 0.07 |
| $\sigma$ | 1.18 | 0.05 |

**Posterior predictive distribution:**

$$\Pr(\log(Inc)|\text{Age} = a)$$

1. Take a sample of size N from posterior $\Pr(\alpha, \beta, \sigma|\text{D})$.
2. For each value of Age $a$, calculate N values of $\mu = \alpha + \beta a$.
3. Draw from Norm($\mu, \sigma$) for each of the N posterior samples.
4. Calculate quantiles of these predicted outcomes.

# Mean versus prediction

## Posterior distribution of mean



## Posterior predictive distribution



For any given age, $\mu$ is the "expected" (mean) log income for people of that age.

The posterior distribution of $\mu$ describes our modeled uncertainty about the value of $\mu$.

This distribution takes into account coefficients $\alpha$ and $\beta$, but not the standard deviation $\sigma$.
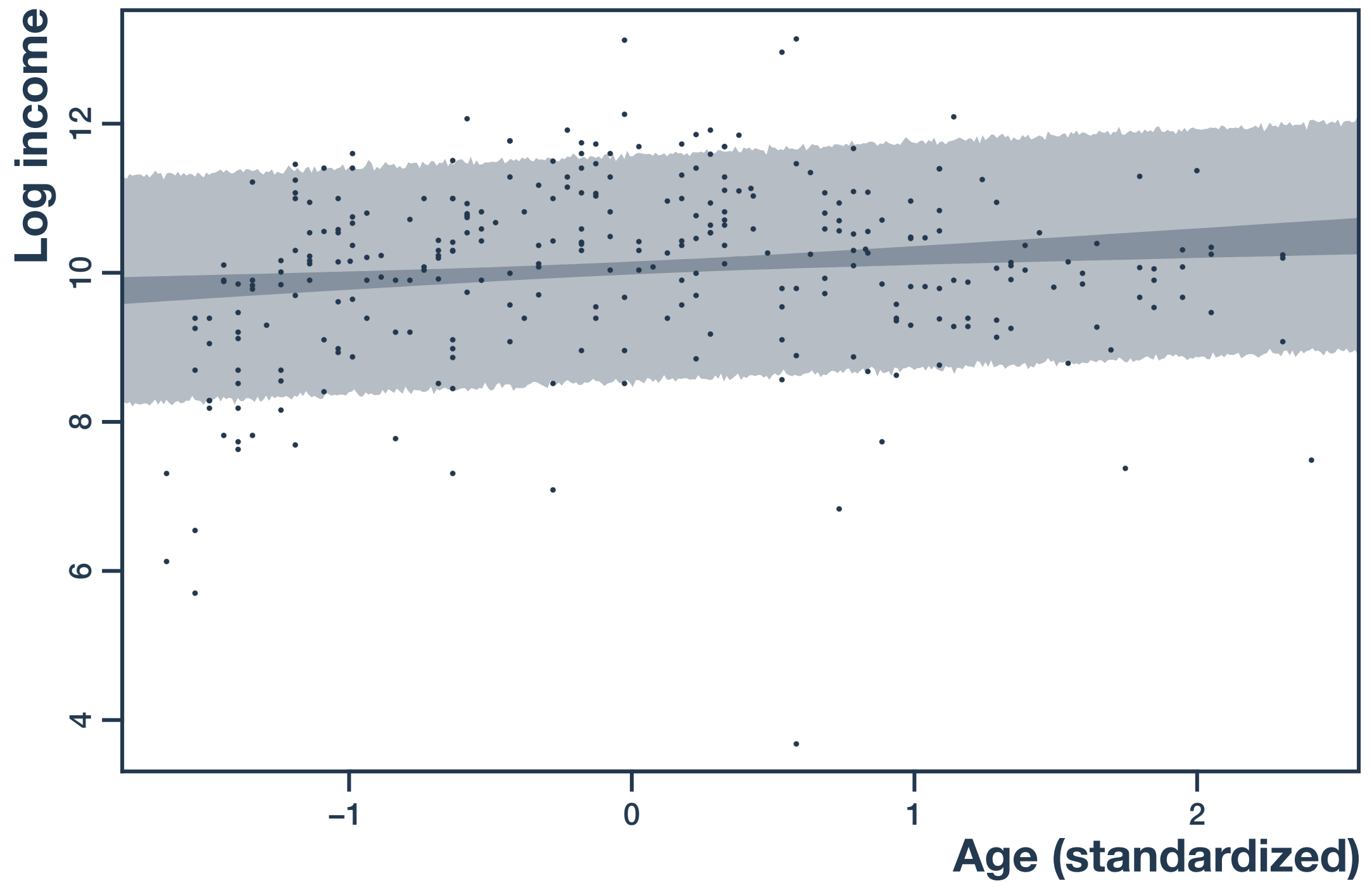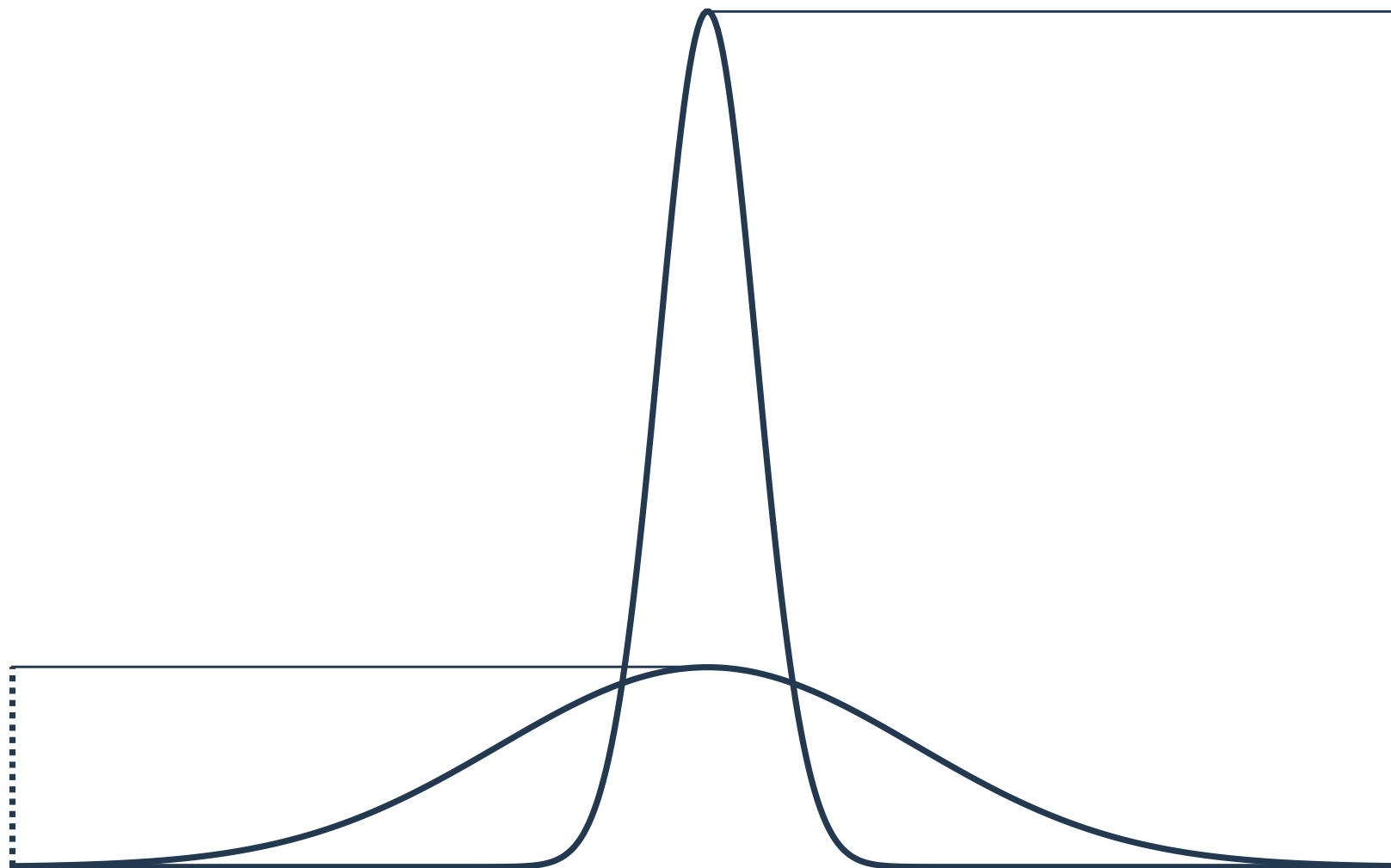
For any given age, the posterior predictive distribution predicts the log income for any individual person of that age.

The posterior distribution of $\mu$ describes our modeled uncertainty about the value of log(Inc).

This distribution takes into account coefficients $\alpha$, $\beta$, and $\sigma$.

The 80% posterior interval should contain about 80% of the data.

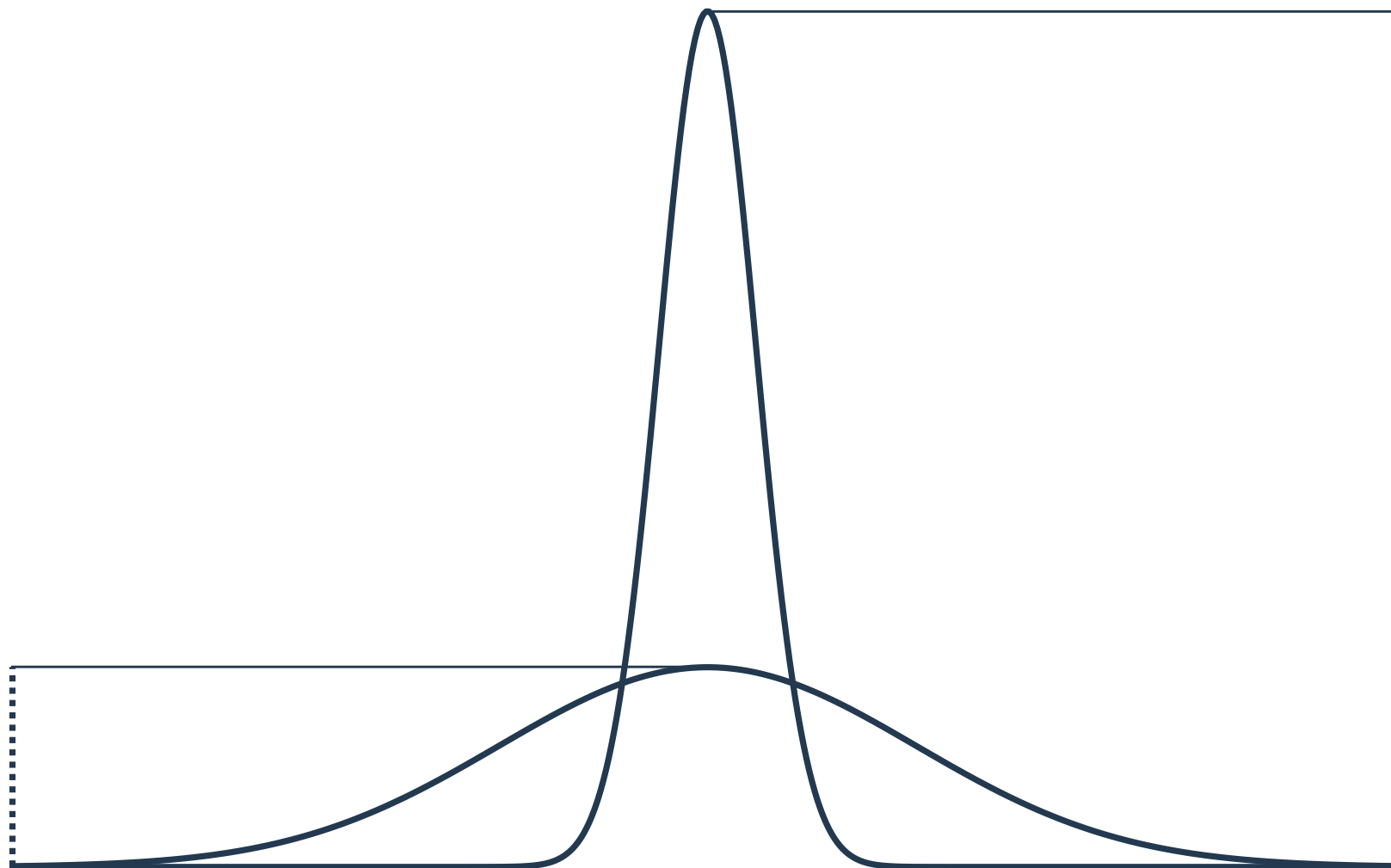# Assessing fit

$$\Pr(\theta|D) = \frac{\Pr(D|\theta)\Pr(\theta)}{\Pr(D)}$$

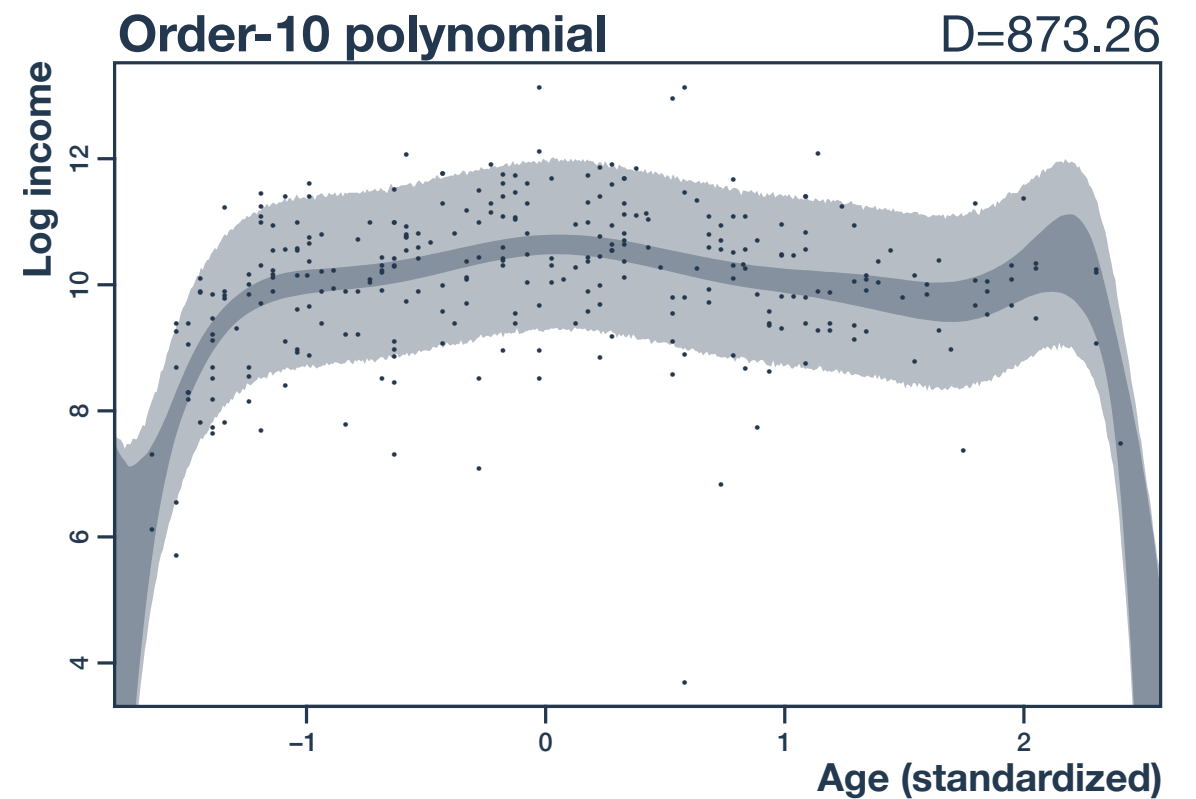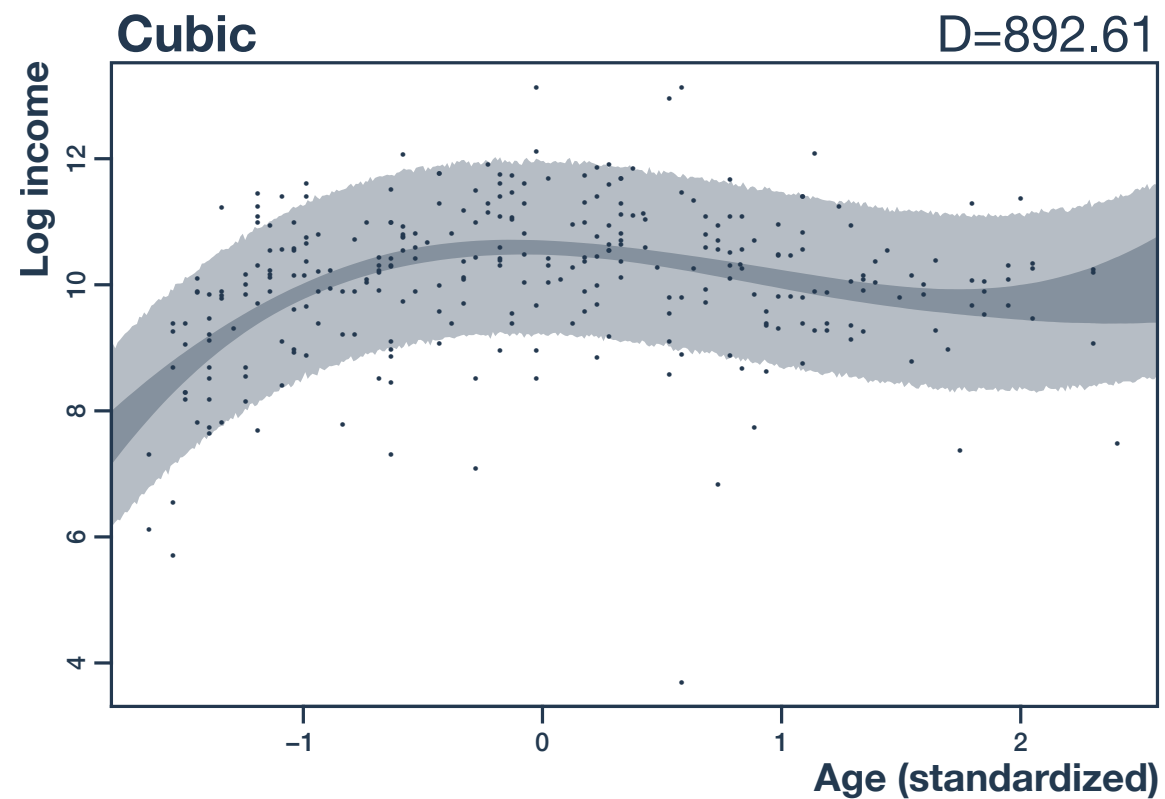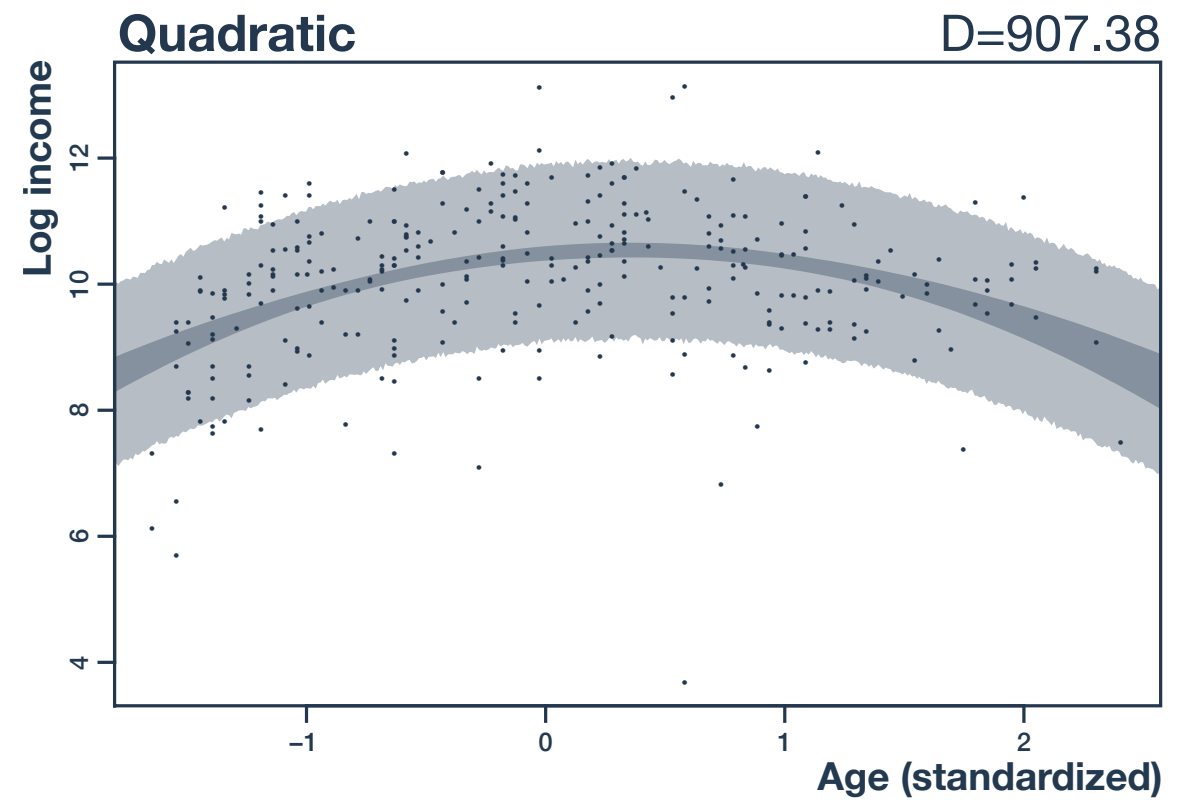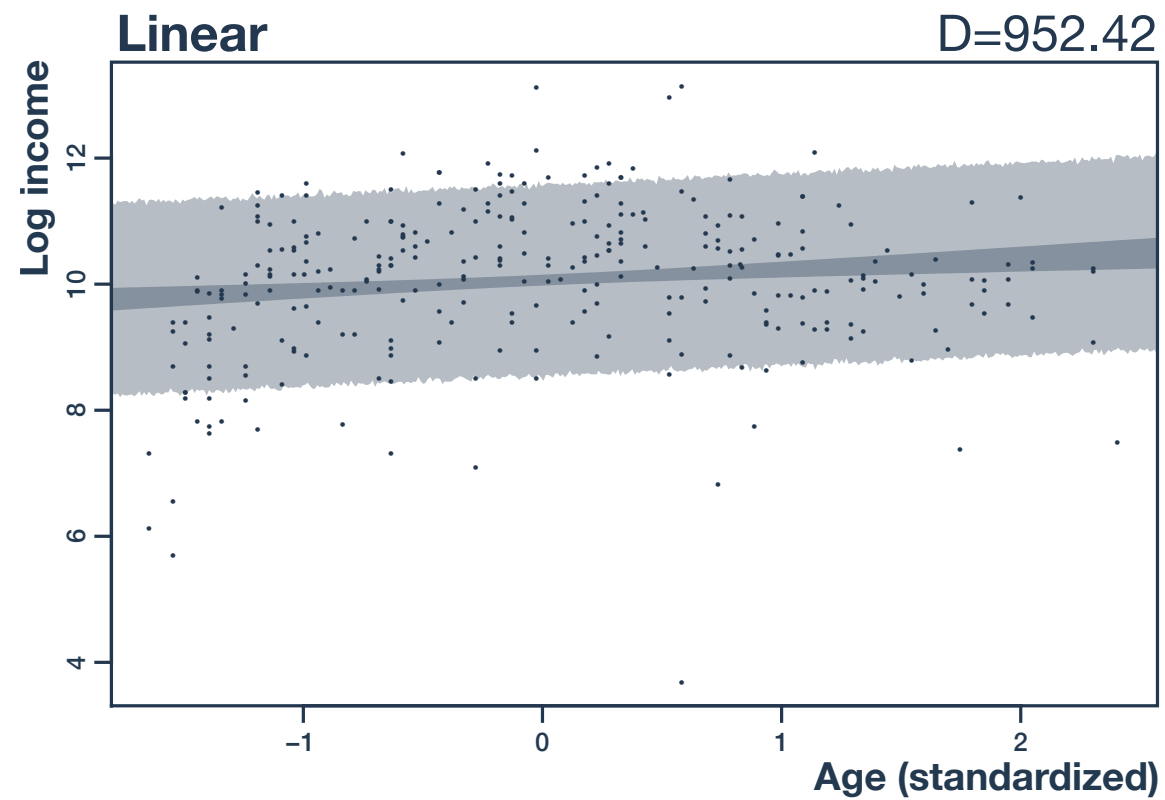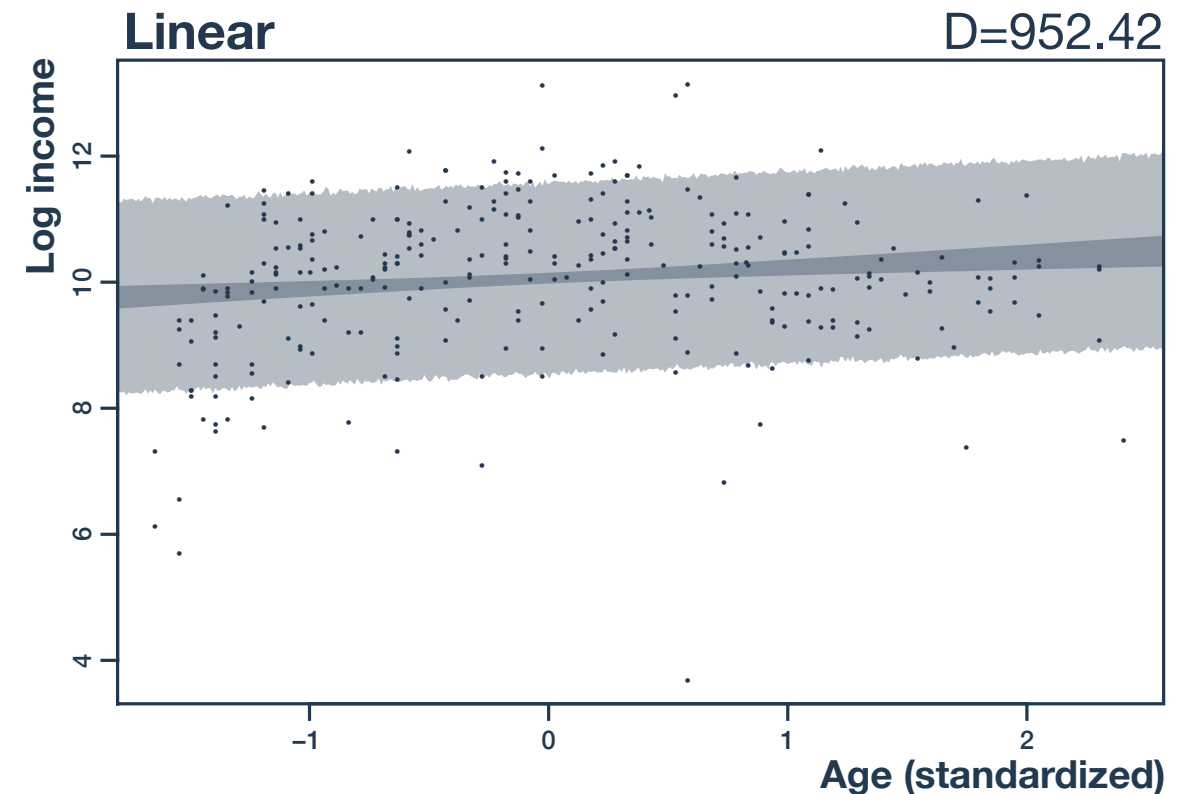$$D = -2\log\big(\Pr(\theta|D)\big)$$

# Deviance

# Goodness of fit

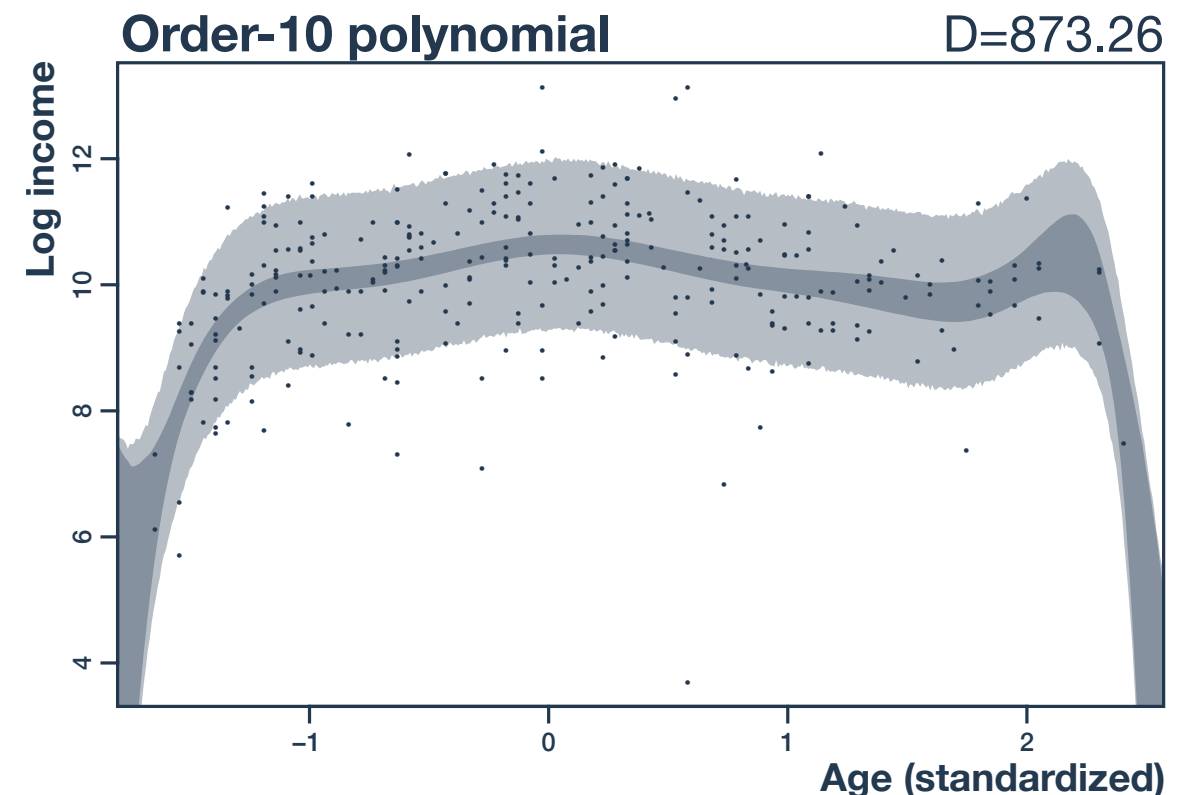**Underfit** | Errs in prediction in a systematic way

Misses important aspects of relationship between predictor and outcome



**Overfit** | Takes random variation to be systematic

Predicts data from sample well, but tends to predict new data very poorly

# Overfitting