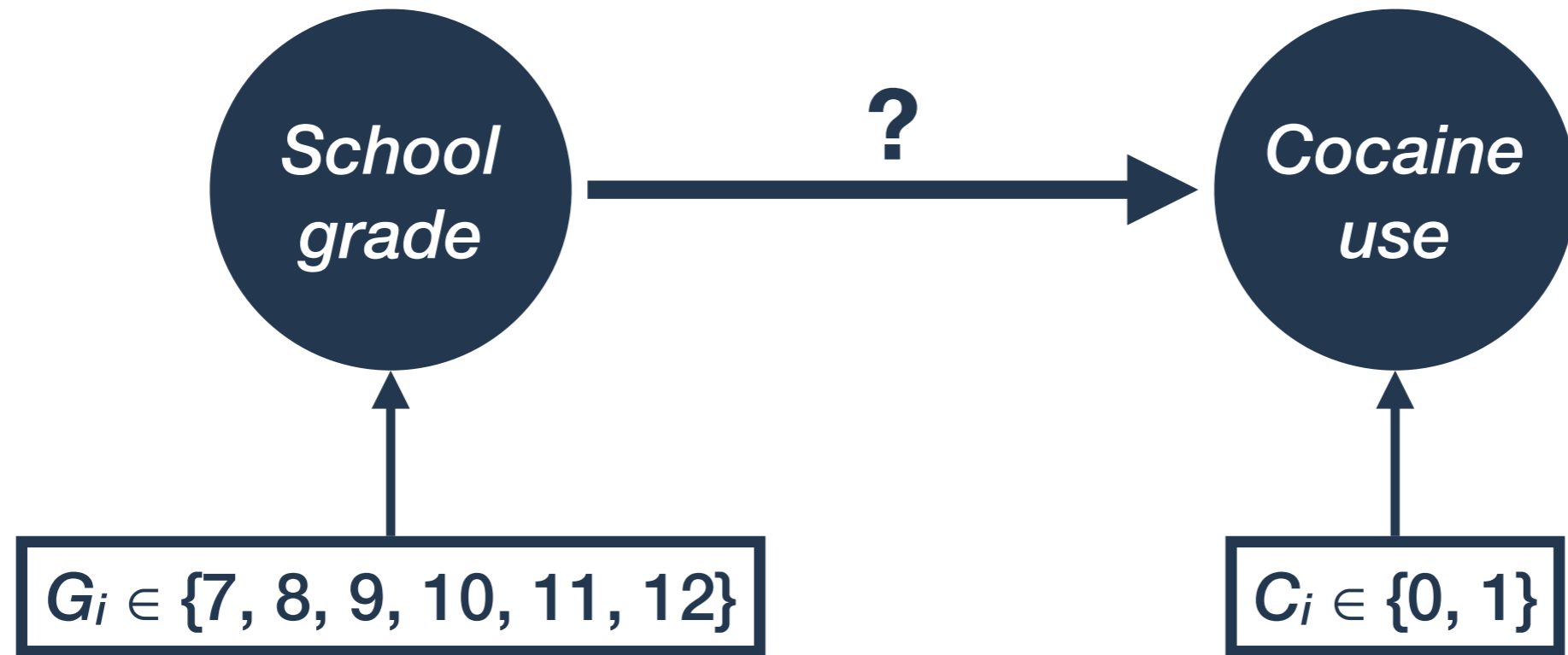


**Feb 2**

Logistic regression:  
motivation

- 1. The trouble with binary outcomes**
- 2. Binomial and Bernoulli distributions**
- 3. Logistic link function**
- 4. Intercept-only logistic regression**

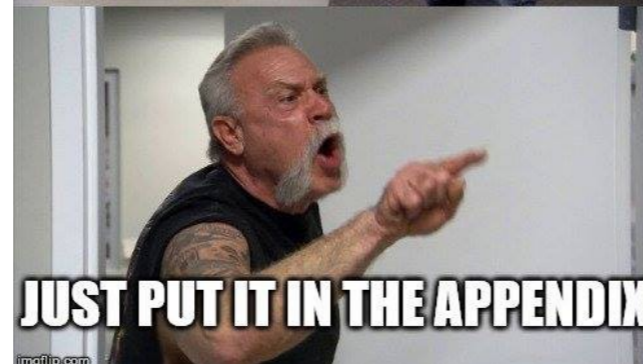
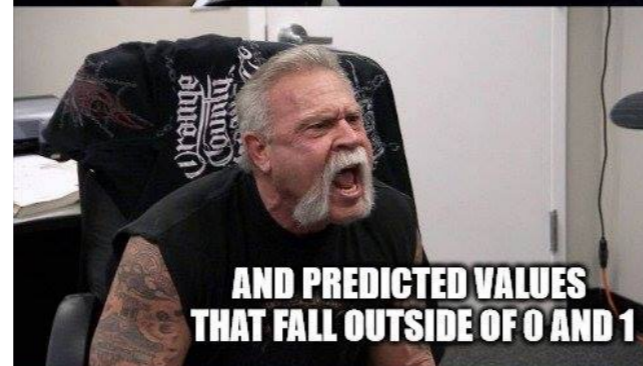
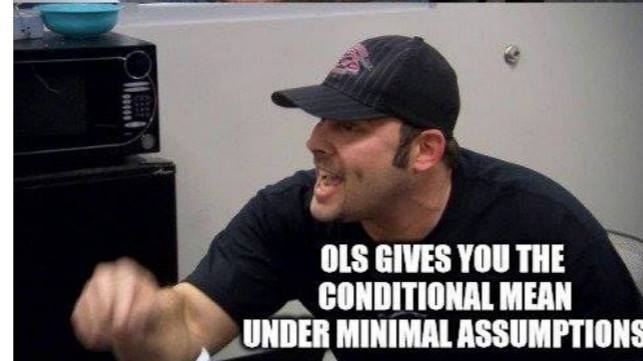
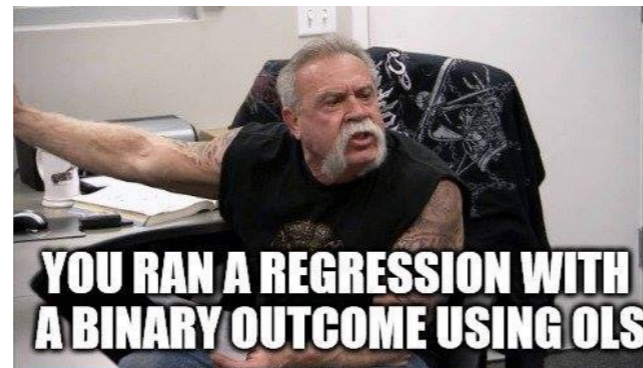
# Cocaine use among adolescents



Why not use a standard linear regression?  $\left| \begin{array}{l} C_i \sim \text{Norm}(\mu, \sigma) \\ \mu = a + \beta G_i \end{array} \right.$

# Normal model of binary data

Why not use a standard linear regression?



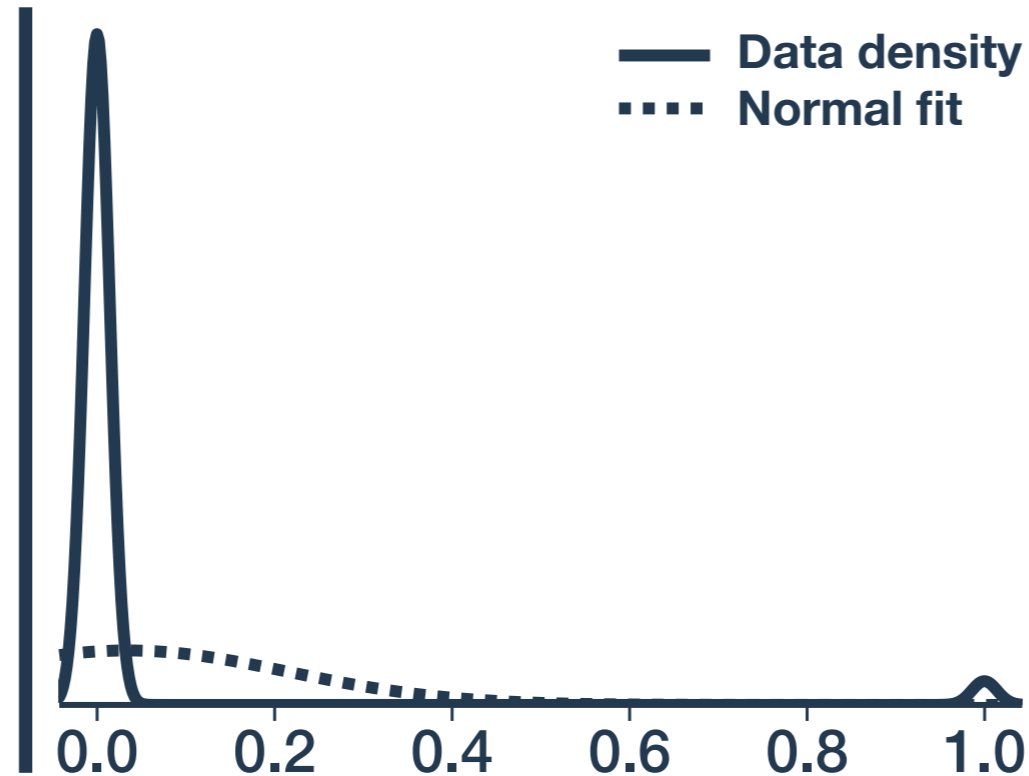
# Normal model of binary data

## Why not use a standard linear regression?

### Wrong support

Normal distribution has a support of  $(-\infty, \infty)$ , but we know the outcome variable takes on only two values.

### Bad intuitive “fit”



### Interpretation

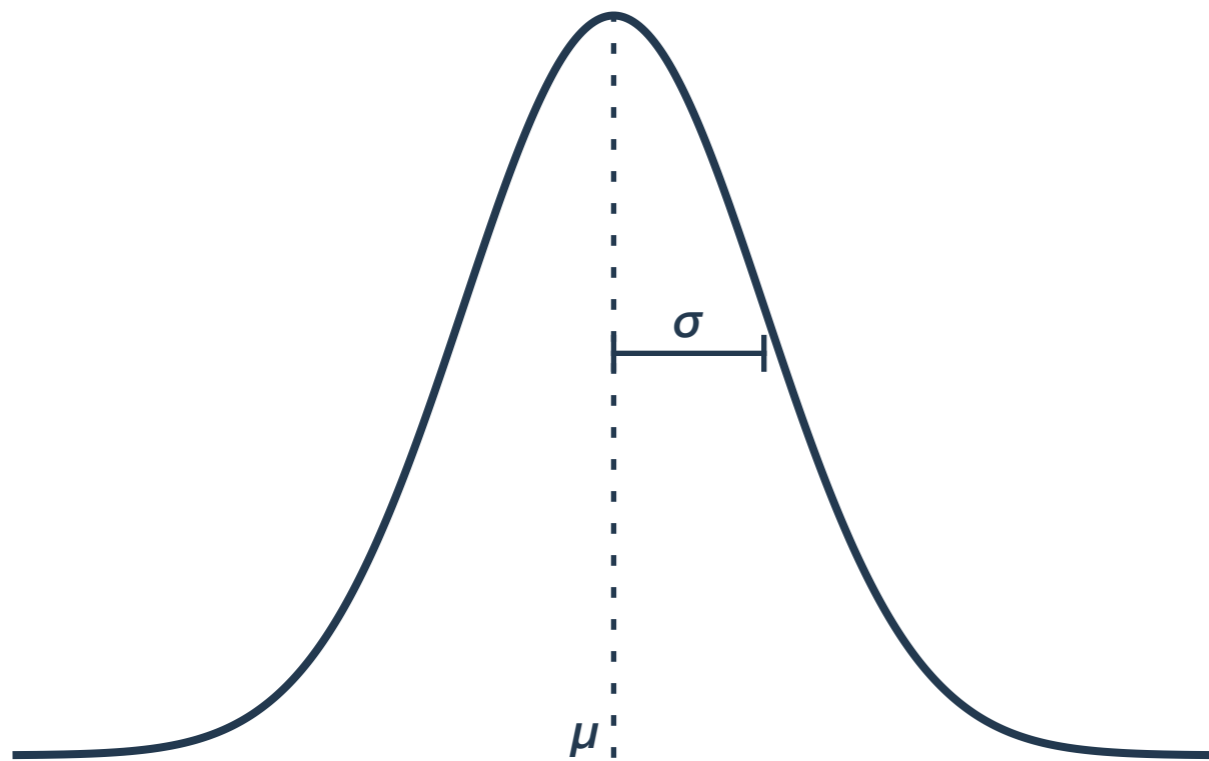
Under some circumstances, results can be interpreted as proportions or probabilities, but this can lead to predicted values less than zero or more than one.

# Normal versus Bernoulli

## Normal distribution

Norm( $\mu, \sigma$ )

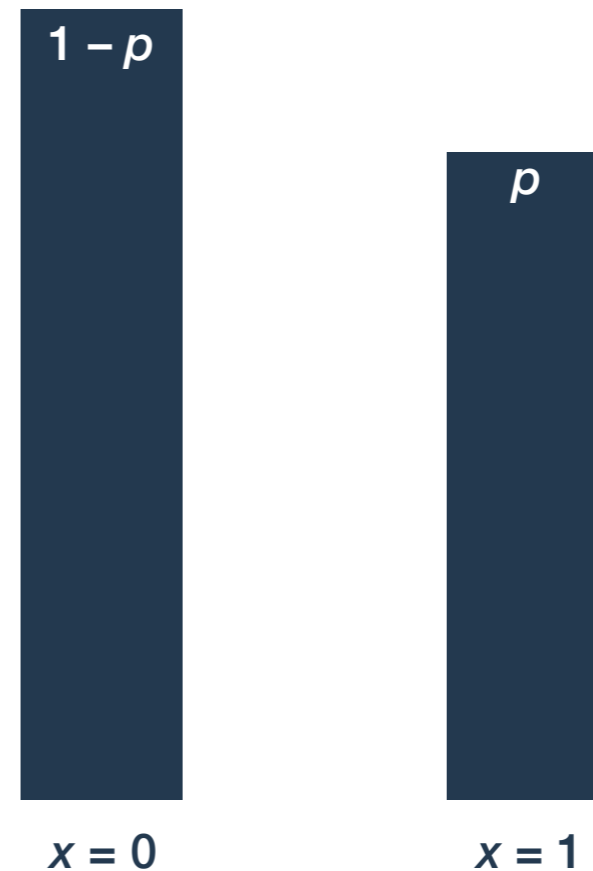
$$\Pr(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



## Bernoulli distribution

Bernoulli( $p$ ) = Binomial(1,  $p$ )

$$\Pr(x|p) = \begin{cases} 1 - p & \text{if } x = 0 \\ p & \text{if } x = 1 \end{cases}$$



# Logistic regression model

Replace  $\text{Norm}(\mu, \sigma)$   
with  $\text{Bernoulli}(p)$



$$C_i \sim \text{Bernoulli}(p_i)$$

$$f(p_i) = a + \beta G_i$$



But now we need  
a “link function”

With normal distribution,  $\mu$   
could take on any value.  
But  $p$  is restricted to  $[0, 1]$ .

# Logistic transformation

## Logit function

$$\text{logit}(p) = \log \left( \frac{p}{1-p} \right)$$

Takes values between 0 and 1, and turns them into values between  $-\infty$  and  $\infty$ .

## Inverse logit function (aka 'logistic')

$$\text{logit}^{-1}(x) = \text{logistic}(x) = \frac{e^x}{e^x + 1} = \frac{1}{1 + e^{-x}}$$

Takes values between  $-\infty$  and  $\infty$ , and turns them into values between 0 and 1.

$$C_i \sim \text{Bernoulli}(p_i)$$

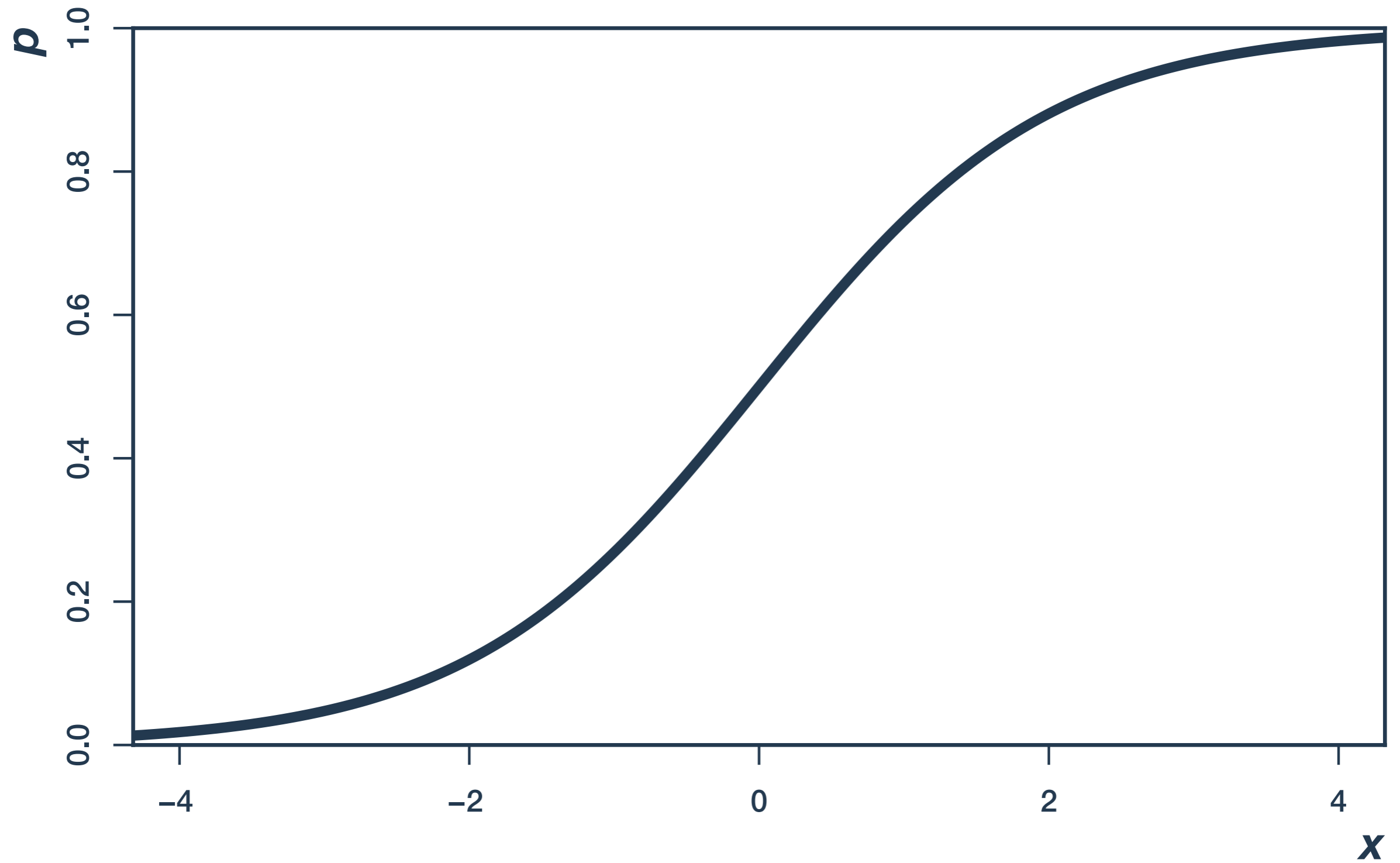
$$\text{logit}(p_i) = a + \beta G_i$$

$$C_i \sim \text{Bernoulli}(p_i)$$

$$p_i = \text{logit}^{-1}(a + \beta G_i)$$

# Logistic transformation

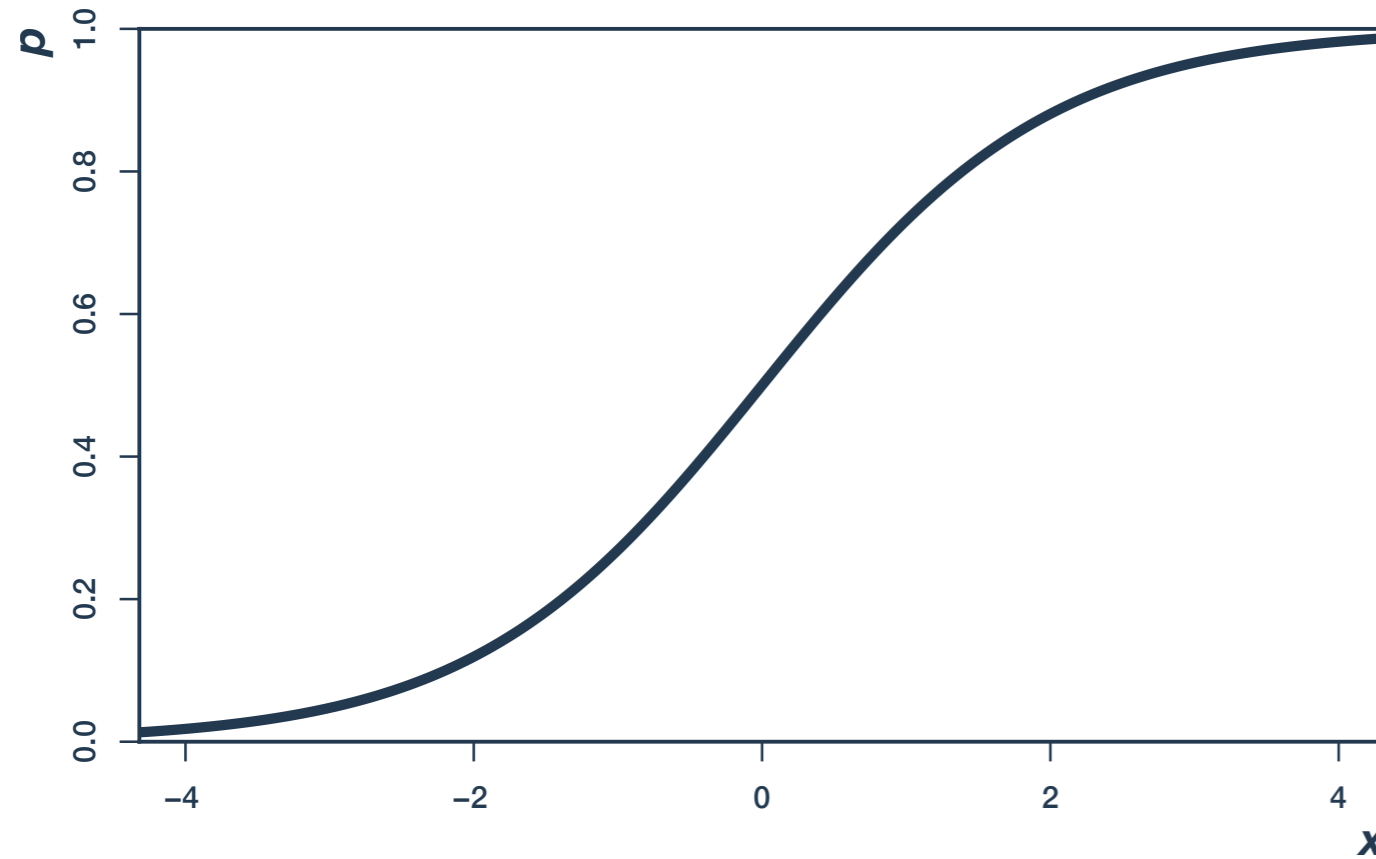
Inverse logit transformation





# Logistic transformation

Inverse logit transformation



$$C_i \sim \text{Bernoulli}(p_i)$$

$$\text{logit}(p_i) = a + \beta G_i$$

---

$x$	$\text{logit}^{-1}(x)$
-2	0.119
-0.5	0.378
0	0.500
0.5	0.622
2	0.881

---

# Intercept-only logistic model

$$C_i \sim \text{Bernoulli}(p_i)$$

$$\text{logit}(p_i) = a$$

Why this model instead of the model we built in the first week of class?

$$\#C \sim \text{Binom}(n, p)$$

$$p \sim \text{Unif}(0, 1)$$

Logistic regression allows us to include explanatory covariates.

# Priors in logistic regression

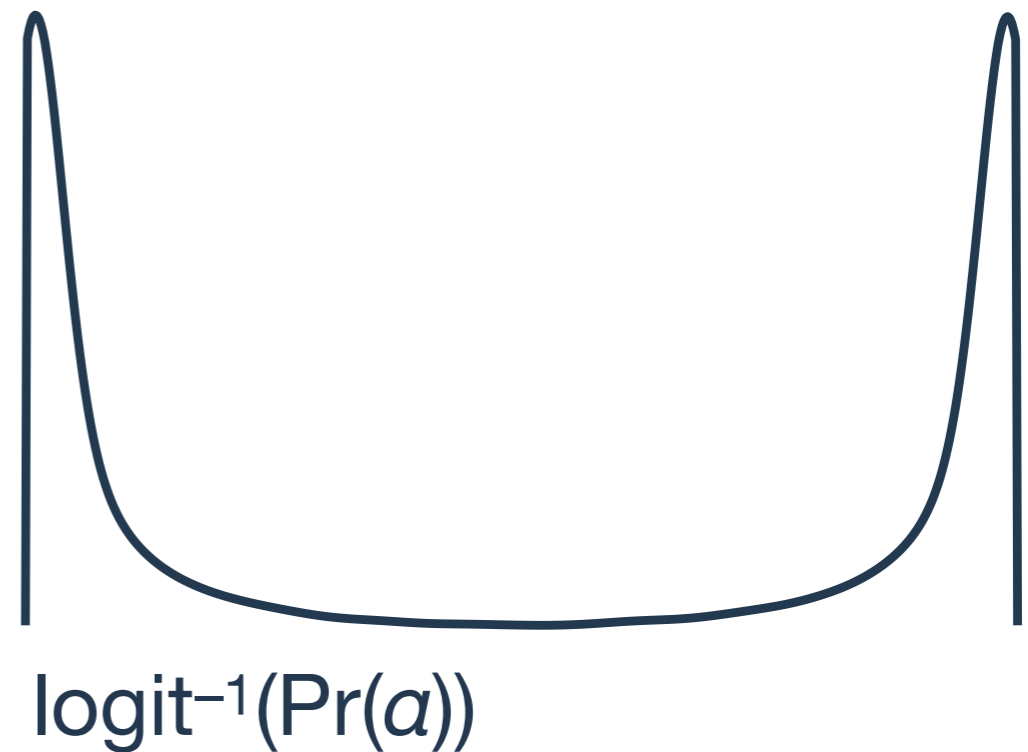
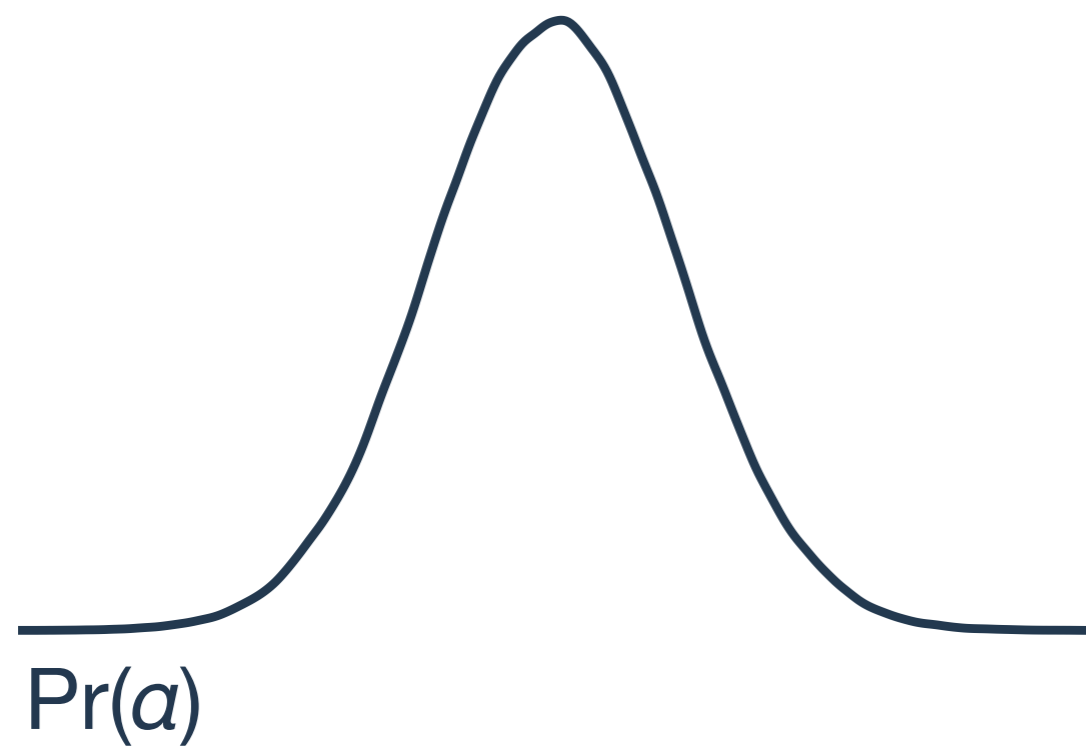
$$C_i \sim \text{Bernoulli}(p_i)$$

$$\text{logit}(p_i) = a$$

$$a \sim \text{Norm}(0, ?)$$

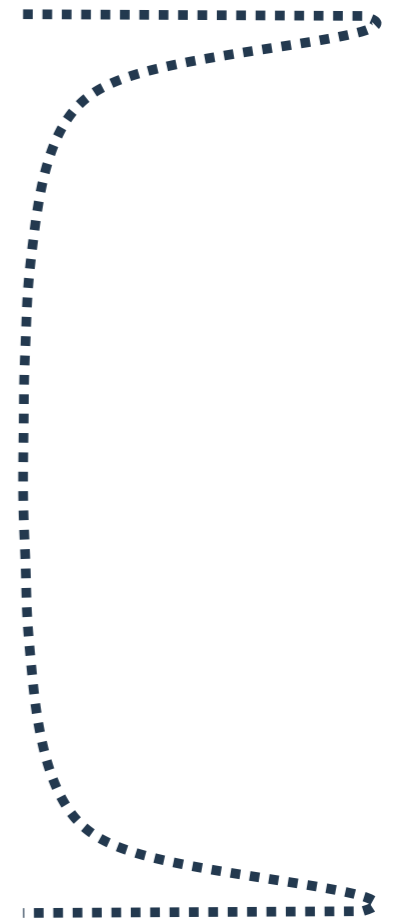
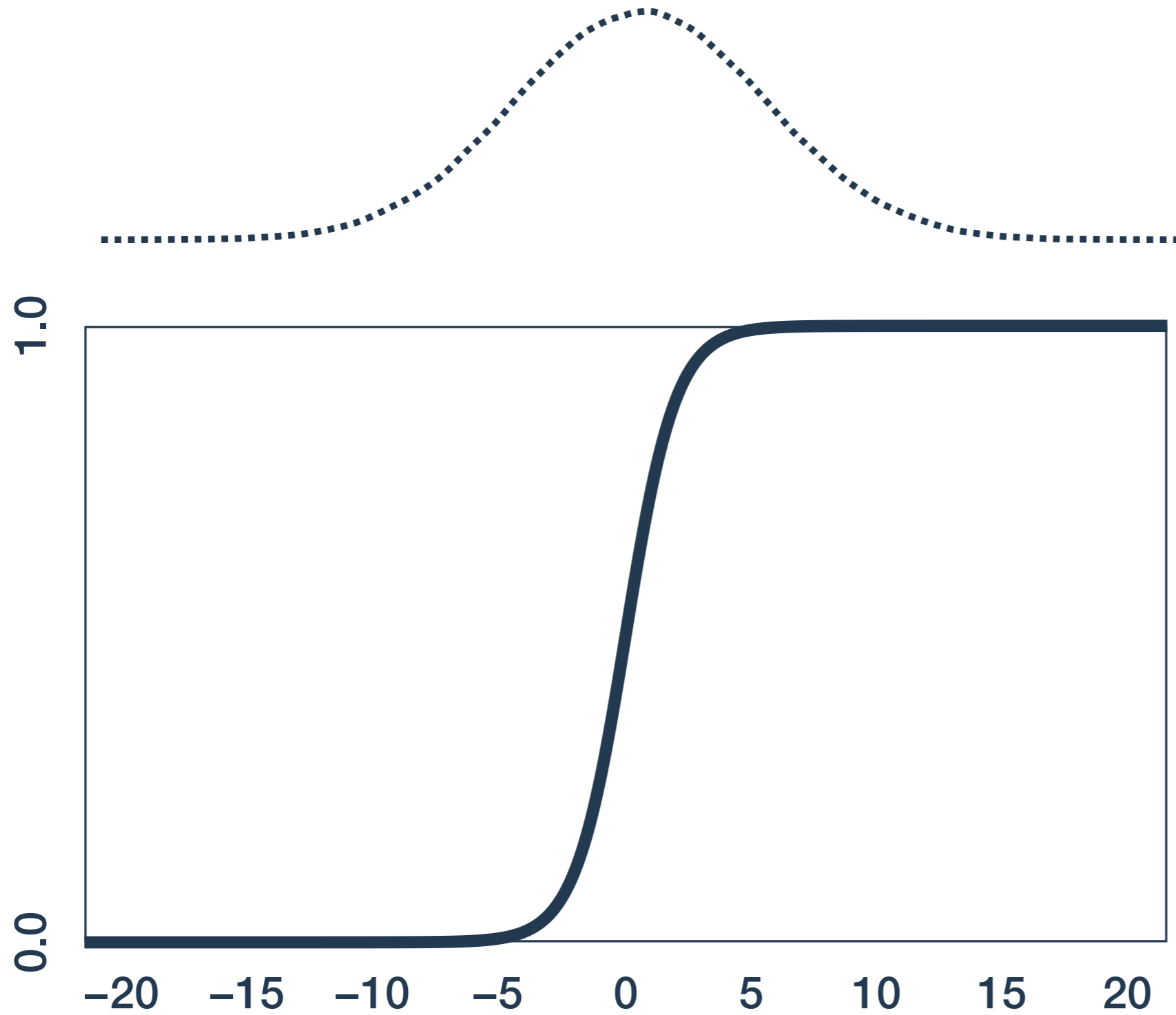
# Priors in logistic regression

$$a \sim \text{Norm}(0, 10)$$



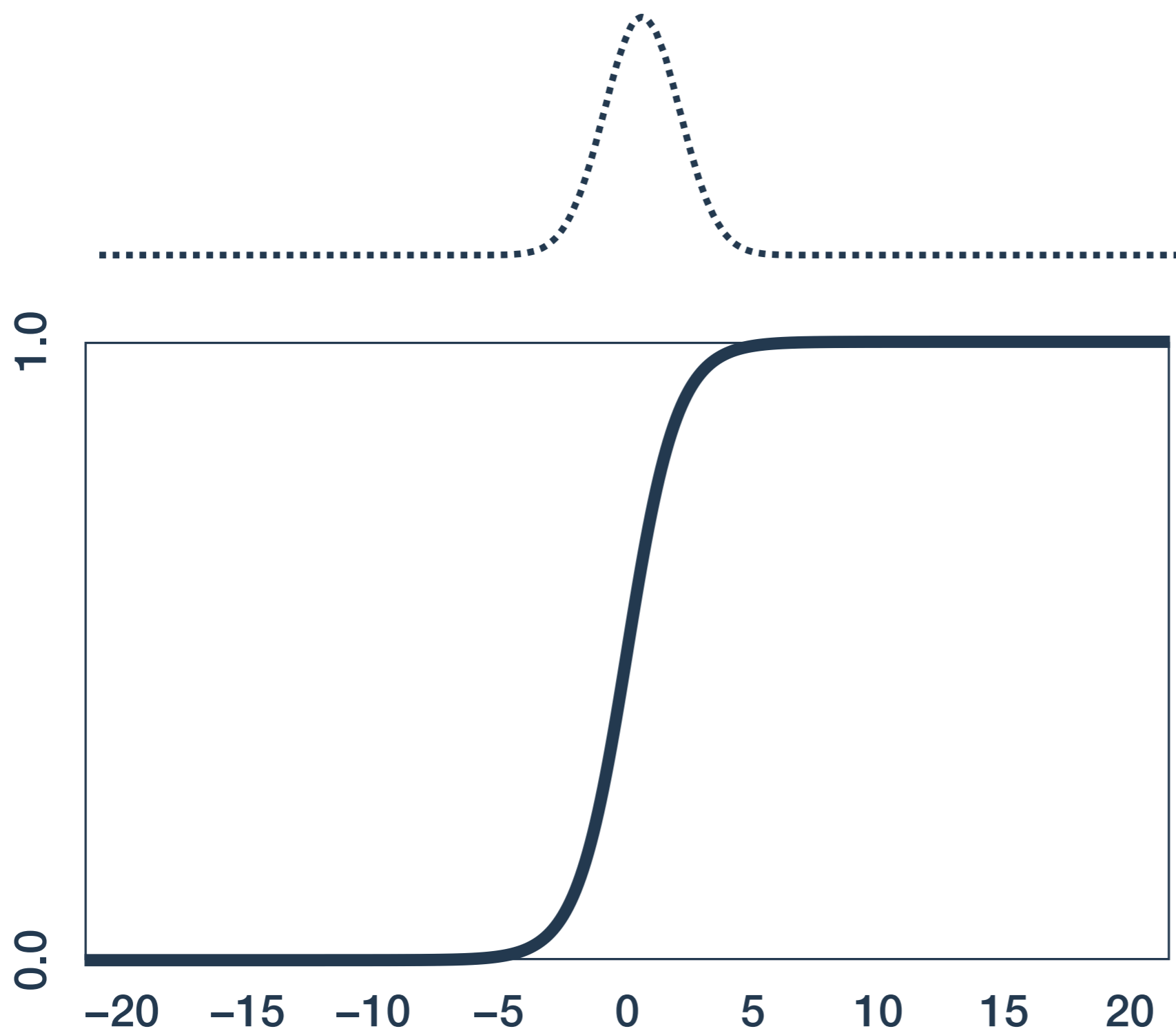
# Priors in logistic regression

$$a \sim \text{Norm}(0, 10)$$



# Priors in logistic regression

$$a \sim \text{Norm}(0, 1.5)$$



# Intercept-only logistic model

$$C_i \sim \text{Bernoulli}(p_i)$$

$$\text{logit}(p_i) = a$$

	Mean	Std. Dev.
<b><i>a</i></b>	-3.334	0.068

n = 6,404

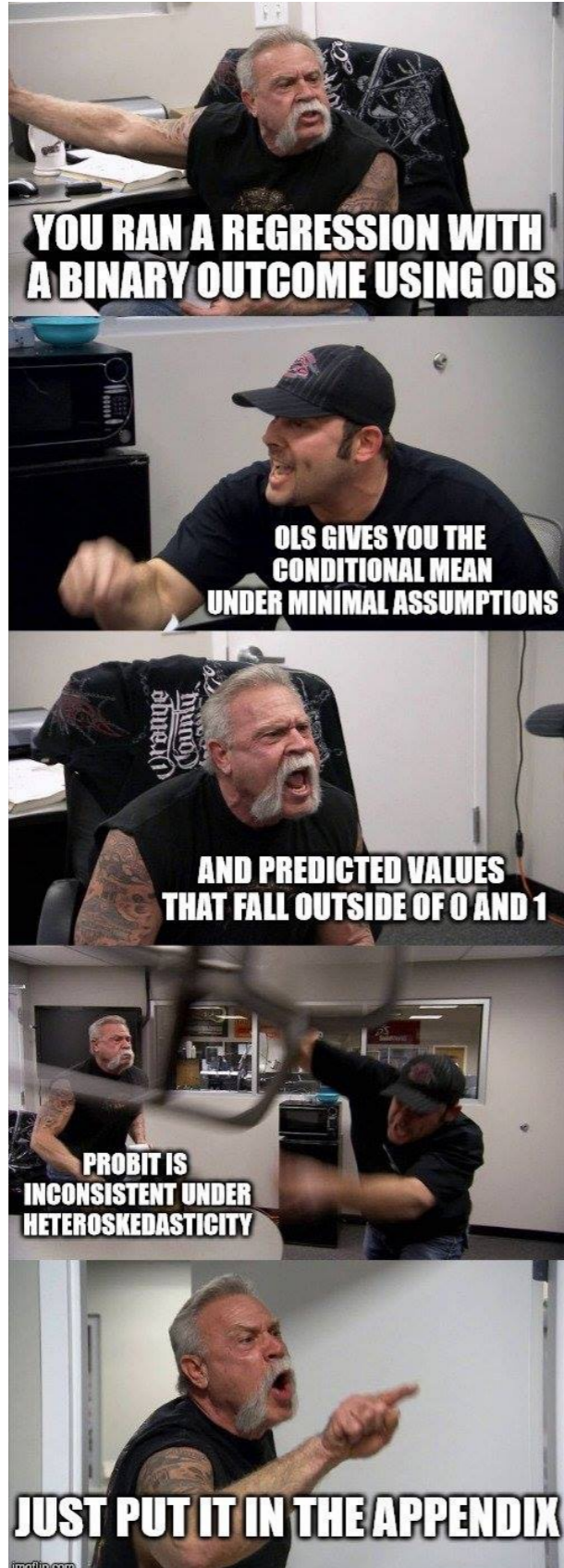
$$a \sim \text{Norm}(0, 1.5)$$

---

$$\exp(-3.334) = 0.0357 \text{ (odds)}$$

$$\text{logit}^{-1}(-3.334) = 0.0344 \text{ (probability)}$$

# Cocaine use among adolescents



Predicted cocaine use among 10th-grade students with standard linear regression

