# Logistic Regression

## BIOE 210

# Binary Classification

Feature $\underline{x}$ → [ Classifier SVM ] → $+1$ / $-1$

## Problems to using a linear model

1. Classification is discrete
2. Bounds

Solution: Logistic Regression

We predict $P(y=1)$, but probabilities $[0,1]$

$$\text{odds}(y) = \frac{P(y=1)}{P(y=0)}$$

$$P(y=1) + P(y=0) = 1 \Rightarrow P(y=0) = 1 - P(y=1)$$

$$\text{odds}(y) = \frac{P(y=1)}{1 - P(y=1)} \quad \in [0, \infty)$$

$$P(y=1) = \frac{\text{odds}(y)}{1 + \text{odds}(y)} \quad \in [0, 1]$$

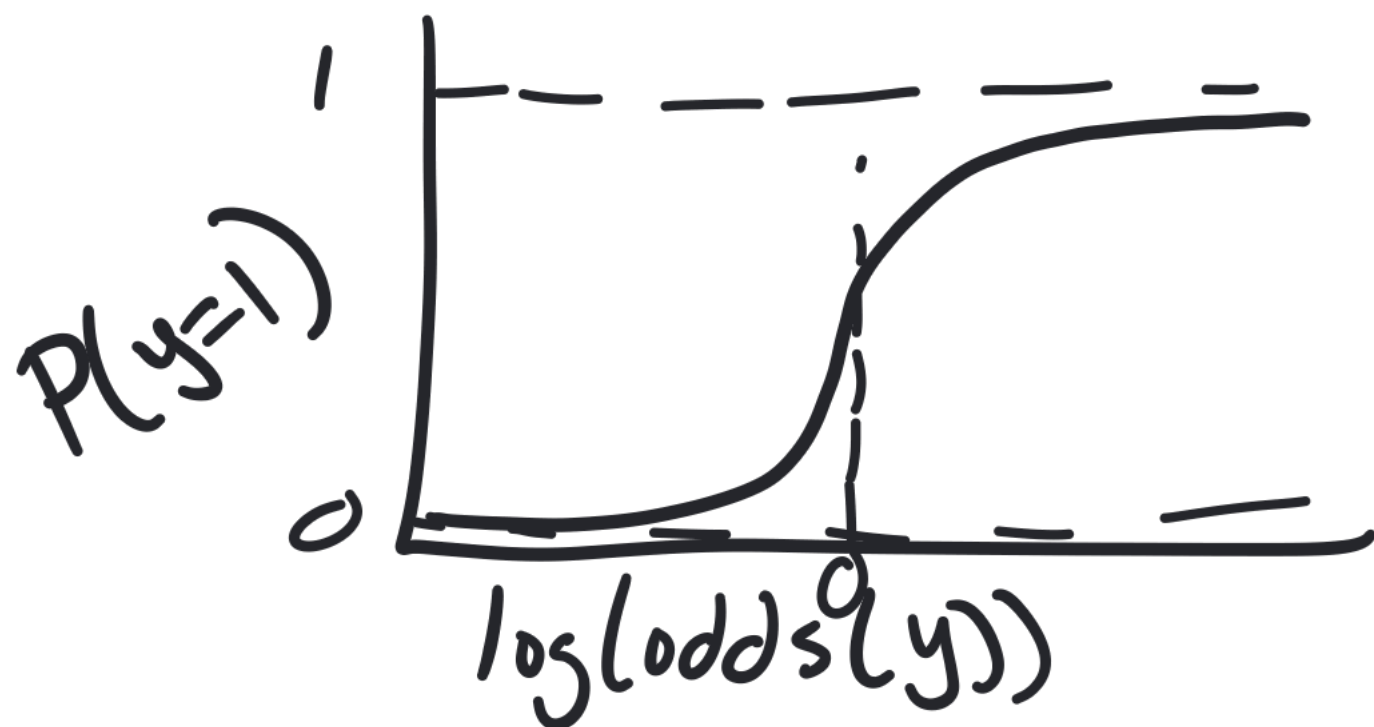$$\log(\text{odds}(y)) \in (-\infty, \infty)$$

# Logistic Regression

$$\log(\text{odds}(y)) = \beta_0 + \beta_1 x_1 + \cdots + \beta_n x_n$$

$$\text{odds}(y) = e^{\beta_0 + \beta_1 x_1 + \cdots + \beta_n x_n} \equiv e^t$$

$$P(y=1) = \frac{\text{odds}(y)}{1 + \text{odds}(y)} = \frac{e^t}{1 + e^t} = \frac{1}{1 + e^{-t}}$$

$$\underline{x} \rightarrow \begin{array}{c} \text{Linear} \\ \text{model} \end{array} \rightarrow t \rightarrow \frac{1}{1 + e^{-t}} \rightarrow P(y=1)$$

$$\beta_0 + \beta_1 x_1 + \cdots$$
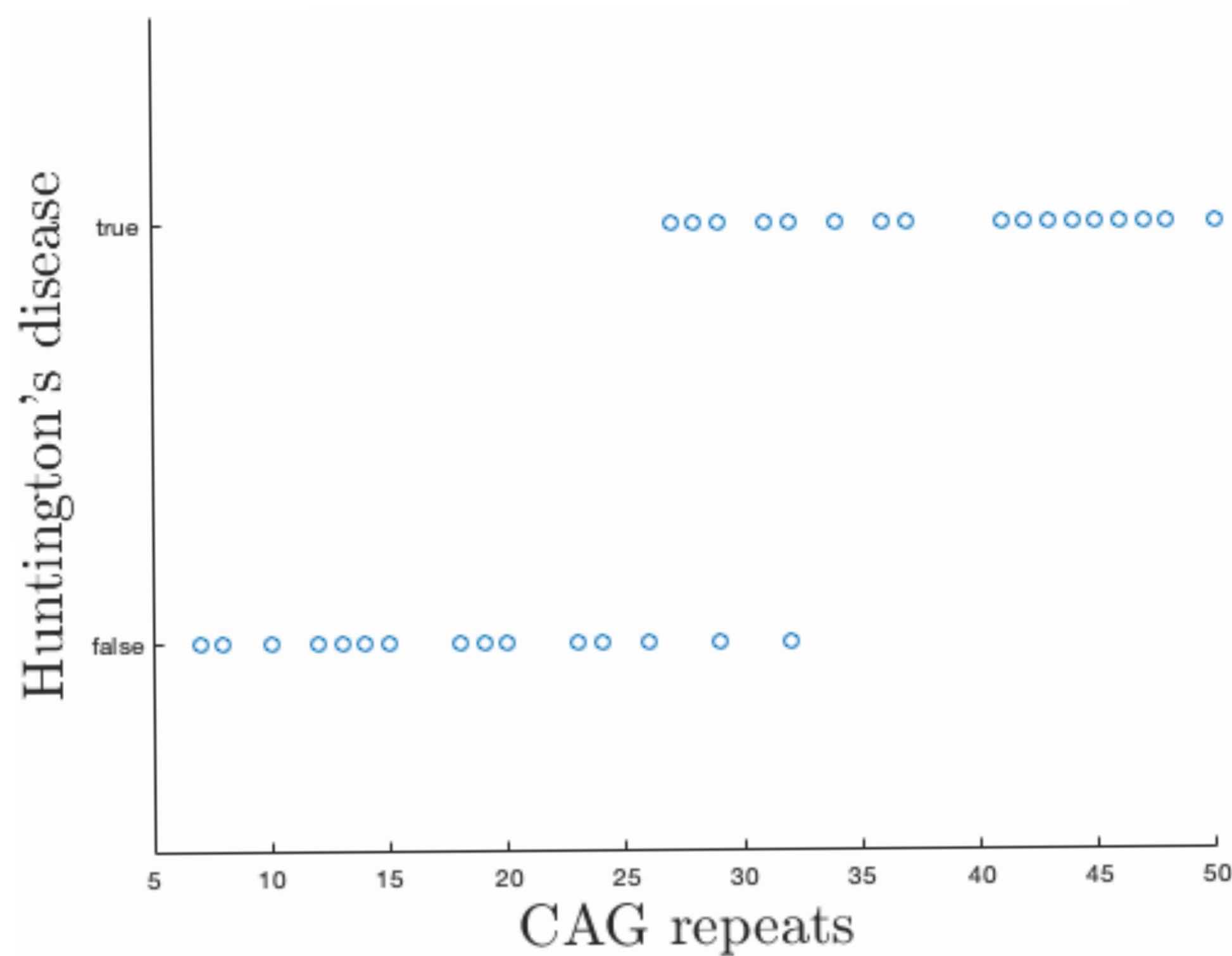
Link function

# Huntingtin (*HTT*)

Leu Lys Ser Phe <u>Gln Gln ... Gln Gln</u> Gln Gln Pro

**ctc aag tcc ttc cag cag ... cag cag caa cag ccg**

| # of CAG Repeats | Disease Outcome |
|---|---|
| < 28 | Not affected. |
| 28-35 | Increased risk. |
| 36-40 | Affected; some offspring affected. |
| > 40 | Affected; all offspring affected. |

Source: Walker FO. Huntington's disease. *The Lancet*. 2007: **369**, (9557), 218–228
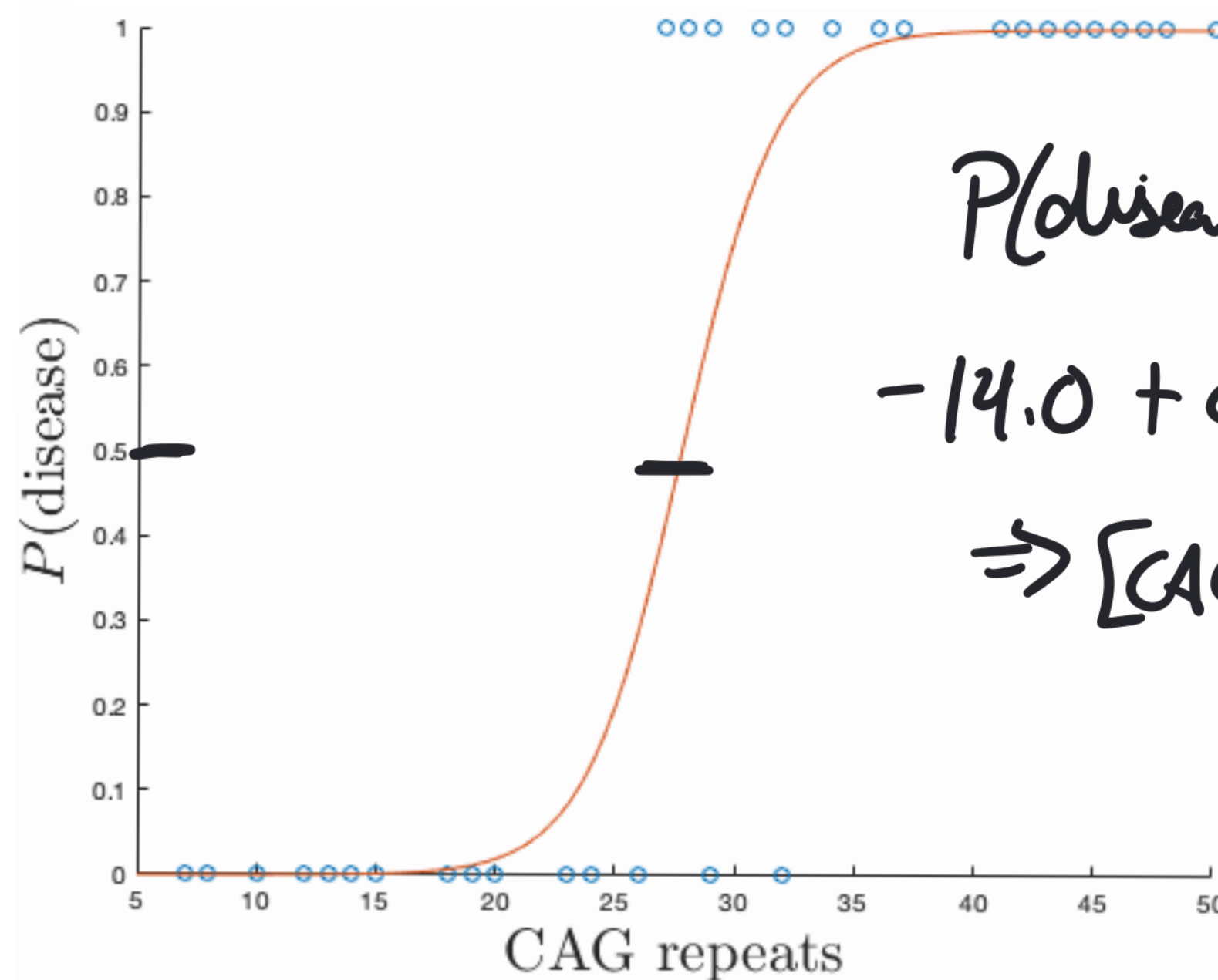
$$P(disease) = f\left([CAGs]\right)$$

$$\log\left(odds(disease)\right) = \beta_0 + \beta_1 [CAGs]$$

$$\log\left(odds(disease)\right) = \beta_0 + \beta_1 [CAGs]$$

$$= -14.0 + 0.51 [CAGs]$$

$$P(disease) = \frac{1}{1 + e^{-14.0 + 0.51 [CAGs]}}$$



$$P(disease) = \frac{1}{2}$$

$$-14.0 + 0.51 [CAGs] = 0$$

$$\Rightarrow [CAGs] = \frac{14.0}{0.51} \approx 28$$

How do we interpret the $\beta$'s?

$$\text{odds ratio}(x_i) = \frac{\text{odds}(x_i + 1)}{\text{odds}(x_i)}$$

$$\text{odds ratio}([CAGs]) = \frac{\text{odds}([CAGs] + 1)}{\text{odds}([CAGs])}$$

$$= \frac{e^{\beta_0 + \beta_1([CAGs] + 1)}}{e^{\beta_0 + \beta_1[CAGs]}}$$

$$= \frac{e^{\beta_0} e^{\beta_1[CAGs]} e^{\beta_1}}{e^{\beta_0} + e^{\beta_1[CAGs]}} = e^{\beta_1}$$

$$\text{odds ratio}(x_i) = \frac{\text{odds}(x_i + 1)}{\text{odds}(x_i)} = e^{\beta_i}$$