# Logistic regression and prediction

# Agenda

- Exam 1
    - Wednesday September 21, in class
    - Covers material up through today (inclusive)
    - Closed notes
    - Bring a calculator (cannot use phone or laptop)
    - I won't ask you to write R code, but you may need to interpret R output
    - Questions similar to assignments and class activities
- Today: more logistic regression

# Data

Data on 5720 Vietnamese children, admitted to hospital with possible dengue fever. Variables include:

- `Dengue`: whether the patient actually has dengue fever, based on a lab test (0 = no, 1 = yes)
- `Temperature`: patient's body temperature (in Celsius)
- `Abdominal`: whether the patient has abdominal pain (0 = no, 1 = yes)
- `HCT`: patient's hematocrit (proportion of red blood cells)
- `Age`: patient's age (in years)
- `Sex`: patient's sex
- + several others

# Last time

$$Y_i \sim Bernoulli(\pi_i)$$

$$\log\left(\frac{\pi_i}{1 - \pi_i}\right) = \beta_0 + \beta_1\, Temperature_i + \beta_2\, Abdominal_i$$

$$+ \beta_3\, Temperature_i \cdot Abdominal_i$$

Does the model improve when we add hematocrit (the proportion of red blood cells)?

# Model

$$Y_i \sim Bernoulli(\pi_i)$$

$$\log\left(\frac{\pi_i}{1 - \pi_i}\right) = \beta_0 + \beta_1\, Temperature_i + \beta_2\, Abdominal_i$$

$$+ \beta_3\, Temperature_i \cdot Abdominal_i$$
$$+ \beta_4\, HCT_i$$

# Class activity, Part I

https://sta214-f22.github.io/class_activities/ca_lecture_11.html

# Class activity

What is the estimated change in odds associated with a 1 point increase in hematocrit, holding temperature and abdominal pain constant?

# Class activity

How does the deviance change when we add hematocrit to the model?

# Class activity

Researchers want to test whether there is a relationship between hematocrit and the probability a patient has dengue, after accounting for temperature and abdominal pain. Carry out a hypothesis test to investigate this research question.

# Comparing models

If deviance always decreases when I add additional variables, how can I assess whether including hematocrit substantially improves the model?

**Option 1:** Likelihood ratio test

✚ Is the change in deviance bigger than we would expect if hematocrit doesn't really matter?

**Option 2:** AIC

# AIC

In linear regression, what quantity did we use to compare models with different numbers of parameters?

# AIC

In linear regression, what quantity did we use to compare models with different numbers of parameters?

*Adjusted $R^2$*

✚ We can use something similar for logistic regression, called the *Akaike information criterion* (AIC)

✚ Motivation: penalize the deviance based on the number of parameters

# AIC

**AIC:** Suppose our model has $p$ parameters (including the intercept). Then the AIC is

$$AIC = 2p + \text{deviance}$$

# AIC

**Model 1:** (adding hematocrit)

```
## Null Deviance:          6956
## Residual Deviance: 6745     AIC: 6755
```

**Model 2:** (no hematocrit)

```
## Null Deviance:          6956
## Residual Deviance: 6914     AIC: 6922
```

Which model do we prefer, based on AIC?

# Model comparison

*Does the model improve when we add hematocrit (the proportion of red blood cells)?*

✚ **Likelihood ratio test:** p-value $\approx 0$

✚ **AIC:** AIC is smaller when we add hematocrit

**Conclusion:** We have convincing evidence that adding hematocrit improves the model.

# A new question...

You report your results to the hospital, and they ask a follow-up question:

*How good is your model at predicting whether a patient has dengue?*

# Making predictions

+ For each patient in the data, we calculate $\widehat{\pi}_i$

+ But, we want to decide which patients to treat. So we need to guess whether patient $i$ has dengue $(Y_i = 1)$ or doesn't $(Y_i = 0)$

How can we turn $\widehat{\pi}_i$ into a dengue prediction?

# Confusion matrix

|  |  | Actual | |
|---|---|---|---|
|  |  | $Y = 0$ | $Y = 1$ |
| **Predicted** | $\widehat{Y} = 0$ | 3957 | 1631 |
|  | $\widehat{Y} = 1$ | 66 | 66 |

✚ For 3957 patients, we correctly predicted they did not have dengue

✚ For 66 patients, we correctly predicted they had dengue

✚ For 1631 patients, we incorrectly predicted they did not have dengue

Did we do a good job at predicting?

# Accuracy

|  |  | Actual | |
| --- | --- | --- | --- |
|  |  | $Y = 0$ | $Y = 1$ |
| **Predicted** | $\widehat{Y} = 0$ | 3957 | 1631 |
|  | $\widehat{Y} = 1$ | 66 | 66 |

$$
\begin{aligned}
\mathrm{Accuracy} &= \frac{\text{number of correct predictions}}{\text{number of observations}} \\
&= \frac{3957 + 66}{5720} \\
&= 0.703
\end{aligned}
$$

**We correctly predict dengue status 70% of the time.**

# Class activity, Part II

https://sta214-f22.github.io/class_activities/ca_lecture_11.html

# Class activity

|  |  | Actual | |
| --- | --- | --- | --- |
|  |  | $Y = 0$ | $Y = 1$ |
| **Predicted** | $\widehat{Y} = 0$ | 3957 | 1631 |
|  | $\widehat{Y} = 1$ | 66 | 66 |

Are our predictions better for patients who actually have dengue, or for patients who don't have dengue?

# Class activity

|  |  | Actual | |
| --- | --- | --- | --- |
|  |  | $Y = 0$ | $Y = 1$ |
| **Predicted** | $\widehat{Y} = 0$ | 3990 | 503 |
|  | $\widehat{Y} = 1$ | 33 | 1194 |

What is the accuracy of the rapid test?

# Class activity

Which method would you prefer -- our logistic regression model, or the rapid test?