# Classification I: Logistic Regression
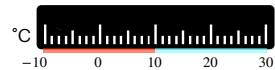Lecture 4

Termeh Shafie

---

## Regression vs. Classification

**continuous variable**
(from $-\infty$ to $\infty$)

$\downarrow$

**categorical variable**
(0 to 1)

- What is the temperature going to be tomorrow?
- Will it be hot or cold tomorrow?

°C

| | | | |
|---|---|---|---|
| $-10$ | 0 | 10 | 20 | 30 |

---

## Linear Probability Model

**predictions**

| | |
|---|---|
| 1 | 0.80 |
| 0 | 0.55 |
| 0 | 0.30 |
| 1 | 0.65 |
| 1 | 0.50 |

**Linear Probability Model vs. Logistic Regression Model**

Predicted Probability of Y = 1

what problems arise when using
linear regression here?

x

---

**Linear Probability Model vs. Logistic Regression Model**

Predicted Probability of Y = 1

x

---

**Logistic Regression: Probability of Heart Disease by Cholesterol Level**
Simulated Data

Probability of Heart Disease

Cholesterol Level

**Probability of passing exam versus hours of studying**

Probability of passing exam (y-axis: 0.00, 0.25, 0.50, 0.75, 1.00)

Hours studying (x-axis: 1, 2, 3, 4, 5)

---

## Redefining The Response

original $Y$

      ●        ●
     0         1

$Y$ as probability

     ●————————●
     0         1

**?**

$Y' \in (-\infty, \infty)$    $-\infty$ ←————————→ $\infty$

how transform $Y$ from $\{0,1\}$ to the real line?

---

## Link Functions

a method to get "non-linear" linear regression
(more on this topic in a later lecture...)

$$y = X\beta \qquad\qquad y = g^{-1}(X\beta)$$

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & x_3 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}$$

the **link function** transforms the probabilities of
the levels of a categorical response variable to
a continuous scale that is unbounded

the **link function** transforms back
the expectation of the response
to the linear function

## Logistic Regression

logit link function and log odds

$$y = X\beta \qquad\qquad y = g^{-1}(X\beta)$$

$$\underbrace{\log\left(\underbrace{\frac{p}{1-p}}_{\text{odds}}\right) = \beta_0 + \beta_1 x_1}_{\text{logit link function}} \qquad p = \frac{e^{\beta_0 + \beta_1 x_1}}{1 + e^{\beta_0 + \beta_1 x_1}} \text{= [a little algebra]}$$
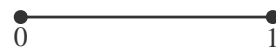
$$= \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1)}}$$

---

## Redefining The Response

original $Y$

$$0 \qquad\qquad 1$$

$Y$ as probability

$$0 \qquad\qquad 1$$

odds of $Y$

$$0 \qquad\qquad \infty$$

$Y' \in (-\infty, \infty)$

$$-\infty \qquad\qquad \infty$$

---

## Logistic Regression

$$y = X\beta \qquad\qquad y = g^{-1}(X\beta) \quad \text{general case}$$

our link function is
$$g(x) = \log \frac{x}{1-x}$$
which has the inverse
$$g^{-1}(x) = \frac{e^x}{1 + e^x}$$

$$p = \frac{e^{X\beta}}{1 + e^{X\beta}} \quad \text{specific case}$$

**Common distributions with typical uses and canonical link functions**

| Distribution | Support of distribution | Typical use | Link name | Link function, $\mathbf{X}\beta = g(\mu)$ | Mean function |
|---|---|---|---|---|---|
| Normal | real: $(-\infty, +\infty)$ | Linear-response data | Identity | $\mathbf{X}\beta = \mu$ | $\mu = \mathbf{X}\beta$ |
| Exponential | real: $(0, +\infty)$ | Exponential-response data, scale parameters | Negative inverse | $\mathbf{X}\beta = -\mu^{-1}$ | $\mu = -(\mathbf{X}\beta)^{-1}$ |
| Gamma | | | | | |
| Inverse Gaussian | real: $(0, +\infty)$ | | Inverse squared | $\mathbf{X}\beta = \mu^{-2}$ | $\mu = (\mathbf{X}\beta)^{-1/2}$ |
| Poisson | integer: $0, 1, 2, \ldots$ | count of occurrences in fixed amount of time/space | Log | $\mathbf{X}\beta = \ln(\mu)$ | $\mu = \exp(\mathbf{X}\beta)$ |
| Bernoulli | integer: $\{0, 1\}$ | outcome of single yes/no occurrence | | $\mathbf{X}\beta = \ln\left(\frac{\mu}{1-\mu}\right)$ | |
| Binomial | integer: $0, 1, \ldots, N$ | count of # of "yes" occurrences out of N yes/no occurrences | | $\mathbf{X}\beta = \ln\left(\frac{\mu}{n-\mu}\right)$ | |
| Categorical | integer: $[0, K)$ K-vector of integer: $[0, 1]$, where exactly one element in the vector has the value 1 | outcome of single K-way occurrence | Logit | | $\mu = \frac{\exp(\mathbf{X}\beta)}{1+\exp(\mathbf{X}\beta)} = \frac{1}{1+\exp(-\mathbf{X}\beta)}$ |
| Multinomial | K-vector of integer: $[0, N]$ | count of occurrences of different types (1, ..., K) out of N total K-way occurrences | | $\mathbf{X}\beta = \ln\left(\frac{\mu}{1-\mu}\right)$ | |

---

# Interpreting Logistic Regression Models

- we want to create a spam filter based on 3921 observations/emails
- simple model, one predictor: `to_multiple`

```
Call:
glm(formula = spam ~ to_multiple, family = binomial, data = email)

Deviance Residuals:
   Min      1Q   Median      3Q     Max
-0.477  -0.477  -0.477  -0.477   2.809

Coefficients:
               Estimate Std. Error z value Pr(>|z|)
(Intercept)    -2.11609    0.05618 -37.665  < 2e-16 ***
to_multipleyes -1.80918    0.29685  -6.095 1.1e-09 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 2437.2  on 3920  degrees of freedom
Residual deviance: 2372.0  on 3919  degrees of freedom
AIC: 2376

Number of Fisher Scoring iterations: 6
```
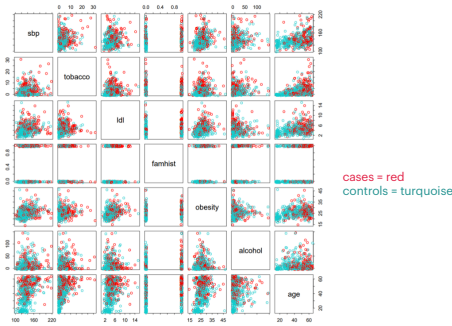
---

# Interpreting Coefficients

| Probability p | Odds p/(1-p) | Log Odds log[p/(1-p)] |
|---|---|---|
| 0.1 | 0.1111 | -2.1972 |
| 0.5 | 1 | 0 |
| 0.9 | 9 | 2.1972 |

## Example: South African Heart Disease

- From Western Cape, South Africa in early 80s
- Coronary Risk Factor Study (CORIS)
- High incidence of myocardial infarction (MI) in region: 5.1%
- Measurements on seven predictors (risk factors)
- 160 cases, 302 controls. Ages 15-64.
- Outcome is presence/absence of MI at time of survey
- Goal:
  - to identify relative strengths and directions of risk factors
  - intervention study aimed at educating the public on healthier diets
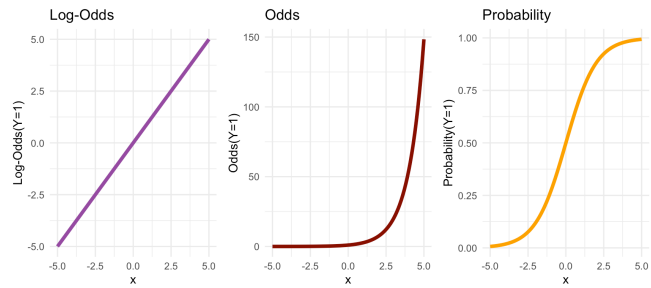
[For more info see ESL 4.4.2]

## Example: South African Heart Disease



cases = red
controls = turquoise

## Example: South African Heart Disease

| term | estimate | std.error | statistic | p.value |
|---|---|---|---|---|
| (Intercept) | -4.130 | 0.964 | -4.283 | 0.000 |
| sbp | 0.006 | 0.006 | 1.023 | 0.306 |
| tobacco | 0.080 | 0.026 | 3.034 | 0.002 |
| ldl | 0.185 | 0.057 | 3.219 | 0.001 |
| famhistPresent | 0.939 | 0.225 | 4.177 | 0.000 |
| obesity | -0.035 | 0.029 | -1.187 | 0.235 |
| alcohol | 0.001 | 0.004 | 0.136 | 0.892 |
| age | 0.043 | 0.010 | 4.181 | 0.000 |

## Interpreting Coefficients
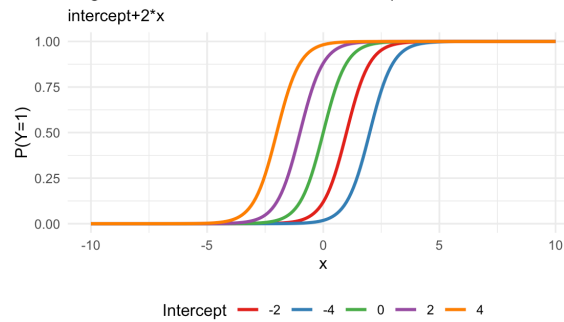
Log-Odds

Odds

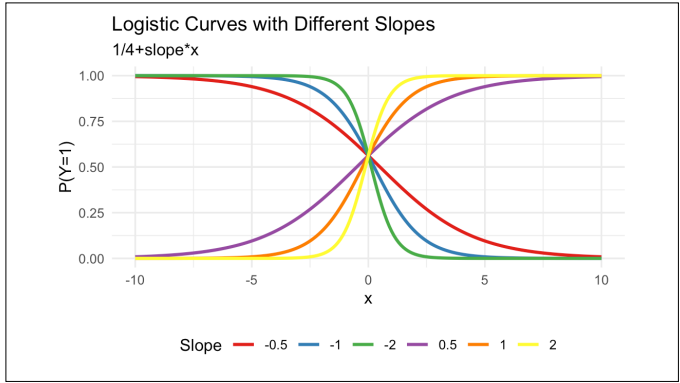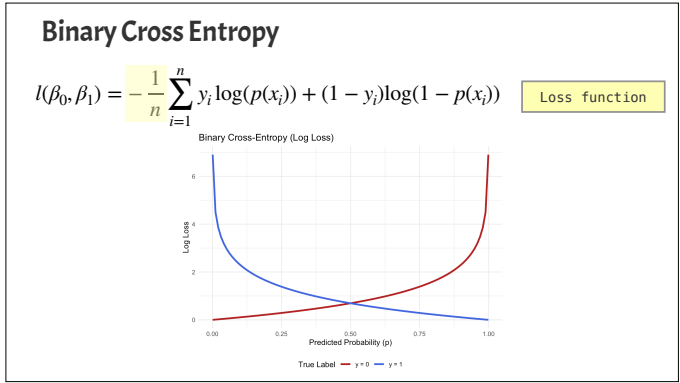Probability

## Estimating Coefficients: MLE

$$\prod_{i;y_i=1} p(x_i) \qquad \prod_{i;y_i=0} 1 - p(x_i)$$

$$L(\beta_0, \beta_1) = \prod_{i;y_i=1} p(x_i) \cdot \prod_{i;y_i=0} 1 - p(x_i)$$

$$L(\beta_0, \beta_1) = \prod_{i=1} p(x_i)^{y_i}(1 - p(x_i))^{1-y_i}$$

$$l(\beta_0, \beta_1) = \sum_{i=1} y_i \log(p(x_i)) + (1 - y_i)\log(1 - p(x_i))$$

[full proof:https://arunaddagatla.medium.com/maximum-likelihood-estimation-in-logistic-regression-f86ff1627b67

## Logistic Curves with Different Intercepts

intercept+2*x

Intercept — -2 — -4 — 0 — 2 — 4

## Logistic Curves with Different Slopes

1/4+slope*x



Slope — -0.5 — -1 — -2 — 0.5 — 1 — 2

## Binary Cross Entropy

$$l(\beta_0, \beta_1) = -\frac{1}{n}\sum_{i=1}^{n} y_i \log(p(x_i)) + (1 - y_i)\log(1 - p(x_i))$$

Loss function



Binary Cross-Entropy (Log Loss)

True Label — y = 0 — y = 1

## Assessing Model Performance

- Did it make the correct prediction?
  - accuracy
  - sensitivity
  - specificity
- How well does to perform in distinguishing classes correctly?

## Confusion Matrix

| | | Predicted | |
|---|---|---|---|
| | | Positive | Negative |
| Actual | Positive | True Positive (TP) | False Negative (TN) |
| | Negative | False Positive (FP) | True Negative (TN) |

## Confusion Matrix: Accuracy

correct predictions

True Positive (TP) + True Negative (TN)

all predictions

False Negative (TN) + False Positive (FP) +

True Positive (TP) + True Negative (TN)

*How often is the model correct?*

| | | Predicted | |
|---|---|---|---|
| | | Positive | Negative |
| Actual | Positive | True Positive (TP) | False Negative (TN) |
| | Negative | False Positive (FP) | True Negative (TN) |

## Confusion Matrix: Sensitivity/Recall

correctly predicted positives
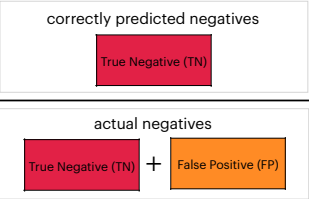
True Positive (TP)

actual positives

False Negative (TN) + True Positive (TP)

*How often is the model correct for Positive Cases?*

| | | Predicted | |
|---|---|---|---|
| | | Positive | Negative |
| Actual | Positive | True Positive (TP) | False Negative (TN) |
| | Negative | False Positive (FP) | True Negative (TN) |

## Confusion Matrix: Specificity

correctly predicted negatives

True Negative (TN)

---

actual negatives

True Negative (TN) + False Positive (FP)

*How often is the model correct for Negative Cases?*

| | | Predicted | |
|---|---|---|---|
| | | Positive | Negative |
| **Actual** | Positive | True Positive (TP) | False Negative (TN) |
| | Negative | False Positive (FP) | True Negative (TN) |

---

## Confusion Matrix: Precision

correctly predicted positives

True Positive (TP)

---

all positives

False Positive (FP) + True Positive (TP)

*How many of the predicted Positives are correct?*

| | | Predicted | |
|---|---|---|---|
| | | Positive | Negative |
| **Actual** | Positive | True Positive (TP) | False Negative (TN) |
| | Negative | False Positive (FP) | True Negative (TN) |

---

## Confusion Matrix: F1 Score

$$\frac{2}{\dfrac{1}{Precision} + \dfrac{1}{Recall}}$$

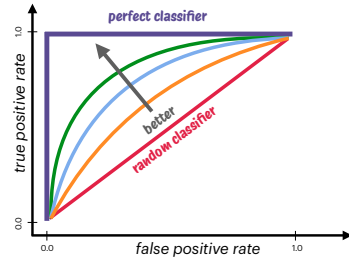*Combination of Precision (how often predicted positives ARE positive) and Recall (how often we correctly predict actual positives)*

$$= \frac{2 \times precision \times recall}{precision + recall}$$

| | | Predicted | |
|---|---|---|---|
| | | Positive | Negative |
| **Actual** | Positive | True Positive (TP) | False Negative (TN) |
| | Negative | False Positive (FP) | True Negative (TN) |

## ROC AUC

correctly predicted positives

True Positive (TP)

actual positives

False Negative (TN) + True Positive (TP)

correctly predicted negatives

True Negative (TN)

$1-$

actual negatives

True Negative (TN) + False Positive (FP)

## ROC AUC



perfect classifier

true positive rate

1.0

0.0

better

random classifier

0.0    false positive rate    1.0

## ROC AUC



true positive rate

false positive rate

Density

Outcome
0
1

Predicted Probability

## ROC AUC



Predicted Probability Density for Banknote Authentication

true positive rate

false positive rate

*more realistic...*

## This Week's Practical

Logistic Regression: The Stock Market Data