Logistics
Splitting your data
Qualitative responses and Logistic Model

# Introduction to Machine Learning for Social Scientists

## Class 5: Intro to GLM

Edgar Franco Vivanco

Stanford University
Department of Political Science

*edgarf1@stanford.edu*

Summer 2018

Logistics
Splitting your data
Qualitative responses and Logistic Model

# Homework 3 Due Friday July 20th at midnight

Available tomorrow Start early!

Logistics
Splitting your data
Qualitative responses and Logistic Model

# In-class midterm Monday July 23rd

Logistics
Splitting your data
Qualitative responses and Logistic Model

# Final project

We'll talk more about this after midterm but you should start forming teams: 6 teams of 5 in total

All teams should be organized by Wed July 25th

Logistics
Splitting your data
Qualitative responses and Logistic Model

# Issues with the Room!

Next Monday old room, this one during the rest of the quarter

Logistics
Splitting your data
Qualitative responses and Logistic Model

# Today's Goals

1. Key concepts
   - Training and test sets
   - Linear Probability Model
   - Logit function and logit inverse function, logistic regression
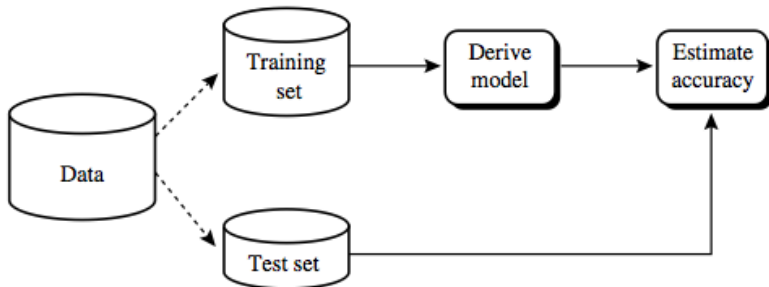2. Key techniques and R functions
   - glm
   - Natural logarithm, log
   - table

Logistics
Splitting your data
Qualitative responses and Logistic Model

# Divide your data

- **Training set:** A set of examples used to fit the parameters and learn. These are already classified by human coders, or produced in a semi-automated way, to create a 'gold standard'.
  Conventionally 80% of available data.

- **Test set:** A set that follows the same probability distribution and is used to test de model.
  Conventionally 20% of available data.

Logistics
Splitting your data
Qualitative responses and Logistic Model

# Divide your data

Logistics
Splitting your data
Qualitative responses and Logistic Model

# R!

Logistics
Splitting your data
Qualitative responses and Logistic Model

# Regression vs Classification

1. *Regression*
    - Quantitative responses
    - Example: Age, height, salary, price, vote share, etc.

Logistics
Splitting your data
Qualitative responses and Logistic Model

# Regression vs Classification

1. *Regression*
   - Quantitative responses
   - Example: Age, height, salary, price, vote share, etc.

2. *Classification*
   - Qualitative responses
   - Example: Election result (win, lose), fake news (yes, no, maybe).

Logistics
Splitting your data
Qualitative responses and Logistic Model

# Regression vs Classification

1. *Regression*
   - Quantitative responses
   - Example: Age, height, salary, price, vote share, etc.
2. *Classification*
   - Qualitative responses
   - Example: Election result (win, lose), fake news (yes, no, maybe).

**NOTE: The distinction is not always that clear-cut:** Logistic regression (a type of non-linear regression) is often used for classification.

Logistics
Splitting your data
Qualitative responses and Logistic Model

# Classification

For simplicity, in this class we are going to classify binary outcomes.

- Vote: Yay / Nay?

Logistics
Splitting your data
Qualitative responses and Logistic Model

# Classification

For simplicity, in this class we are going to classify binary outcomes.

- Vote: Yay / Nay?
- Email: Spam / Not Spam?

Logistics
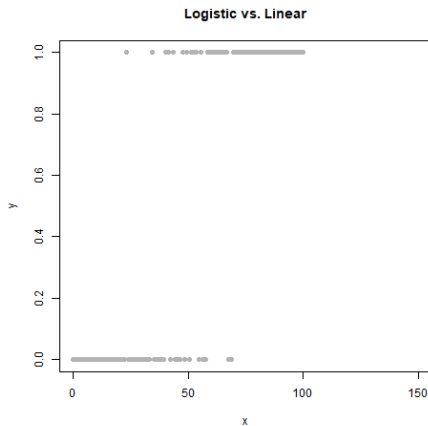Splitting your data
Qualitative responses and Logistic Model

# Classification

For simplicity, in this class we are going to classify binary outcomes.

- Vote: Yay / Nay?

- Email: Spam / Not Spam?

- Online transaction: Faudulent (Yes / No)?

Logistics
Splitting your data
Qualitative responses and Logistic Model

# Classification

For simplicity, in this class we are going to classify binary outcomes.

- Vote: Yay / Nay?

- Email: Spam / Not Spam?

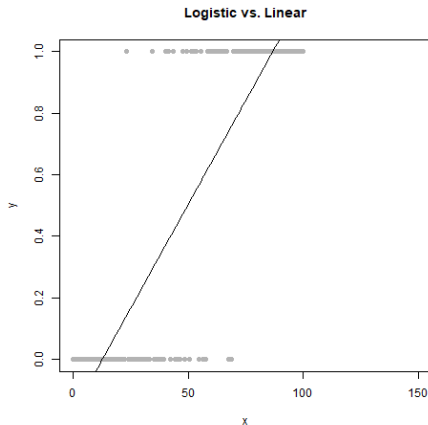- Online transaction: Faudulent (Yes / No)?

$$y \in \{0, 1\} \qquad \left[ \begin{array}{ll} 0: & \text{"Negative class"} \\ 1: & \text{"Positive class"} \end{array} \right]$$
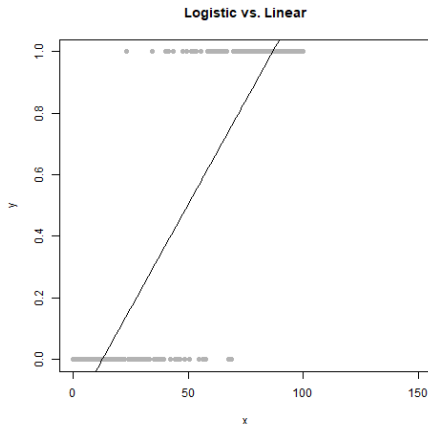
Logistics
Splitting your data
Qualitative responses and Logistic Model

# Predicting qualitative responses



**Logistic vs. Linear**

- $y \in \{0, 1\}$, X continuous

Logistics
Splitting your data
Qualitative responses and Logistic Model

# Predicting qualitative responses



- $y \in \{0, 1\}$, X continuous

- We can always fit a linear model as we did before

Logistics
Splitting your data
Qualitative responses and Logistic Model

# Predicting qualitative responses



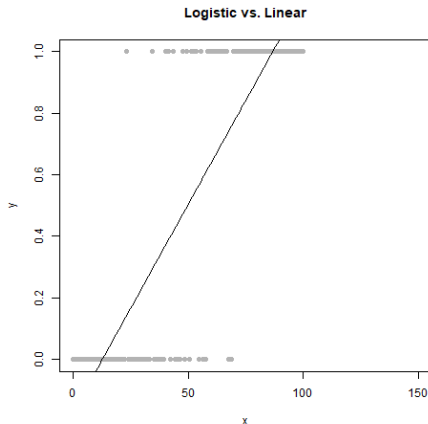- $y \in \{0, 1\}$, X continuous

- We can always fit a linear model as we did before

- We call this the linear probability model because we can associate its predictions to probabilities.

Logistics
Splitting your data
Qualitative responses and Logistic Model

# Predicting qualitative responses



- $y \in \{0, 1\}$, X continuous

- We can always fit a linear model as we did before

- We call this the linear probability model because we can associate its predictions to probabilities.

- PROBLEMS?

Logistics
Splitting your data
Qualitative responses and Logistic Model

# Predicting qualitative responses



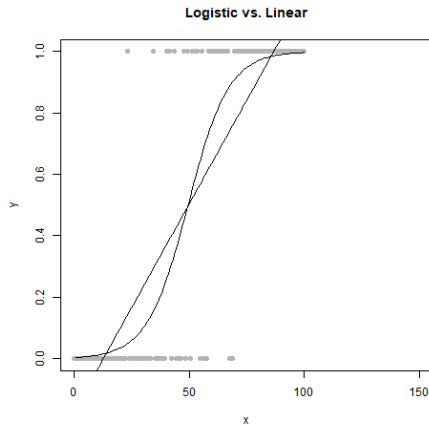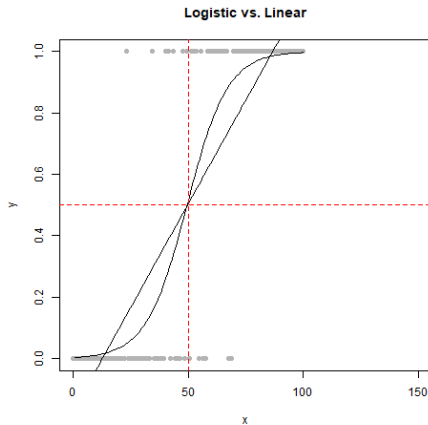- $y \in \{0, 1\}$, X continuous

- We can always fit a linear model as we did before

- We call this the linear probability model because we can associate its predictions to probabilities.

- PROBLEMS?

- Predictions smaller than zero and larger than 1!

Logistics
Splitting your data
Qualitative responses and Logistic Model

# Predicting qualitative responses



Logistic vs. Linear
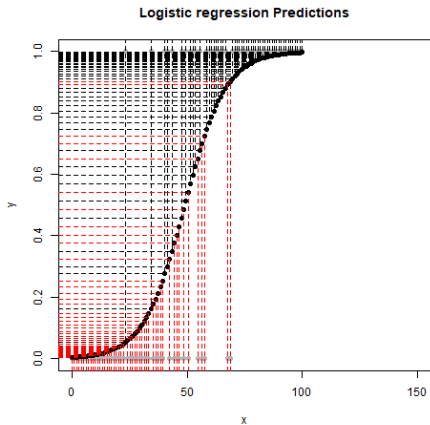
- ▶ Logistic model
- ▶ Estimates the probability of yes: $Pr(\hat{Vote_i} = 1|x_i)$ using a logarithmic transformation

Logistics
Splitting your data
Qualitative responses and Logistic Model

# Predicting qualitative responses



- Logistic model

- Estimates the probability of yes: $Pr(\hat{Vote_i} = 1|x_i)$

- The linear model predictions intersect when the probability equals 0.5.

Logistics
Splitting your data
Qualitative responses and Logistic Model

# Predicting qualitative responses



**Logistic regression Predictions**

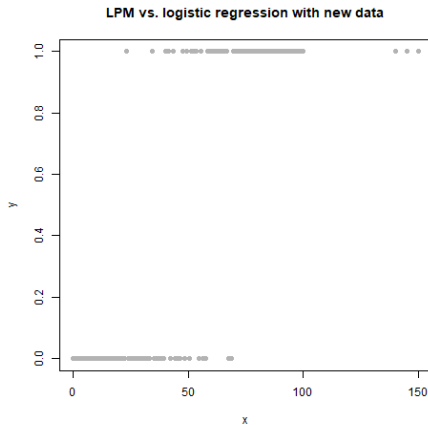- Logistic model
- Estimates the probability of yes: $Pr(\hat{Vote_i} = 1 | x_i)$
- The linear model predictions intersect when the probability equals 0.5.
- We can associate each point to a probability

Logistics
Splitting your data
Qualitative responses and Logistic Model

# Predicting qualitative responses



LPM vs. logistic regression with new data

▶ Logistic model usually more stable to outliers

Logistics
Splitting your data
Qualitative responses and Logistic Model

# Predicting qualitative responses



LPM vs. logistic regression with new data

▶ Logistic model usually more stable to outliers

Logistics
Splitting your data
Qualitative responses and Logistic Model

# A Brief Reminder About (Natural) Logarithms

Logarithm log is a class of functions.

Logistics
Splitting your data
Qualitative responses and Logistic Model

# A Brief Reminder About (Natural) Logarithms

Logarithm log is a class of functions.

- $\log_e z = x$ if $e^x = z$.

Logistics
Splitting your data
Qualitative responses and Logistic Model

# A Brief Reminder About (Natural) Logarithms

Logarithm log is a class of functions.

- $\log_e z = x$ if $e^x = z$.
- We'll call $\log_e$ natural logarithm. And we'll assume $\log = \log_e$

Logistics
Splitting your data
Qualitative responses and Logistic Model

# A Brief Reminder About (Natural) Logarithms

Logarithm log is a class of functions.

- $\log_e z = x$ if $e^x = z$.

- We'll call $\log_e$ natural logarithm. And we'll assume $\log = \log_e$

- $\log e = 1$ (because $e^1 = e$)

Logistics
Splitting your data
Qualitative responses and Logistic Model

# A Brief Reminder About (Natural) Logarithms

Logarithm log is a class of functions.

- $\log_e z = x$ if $e^x = z$.
- We'll call $\log_e$ natural logarithm. And we'll assume $\log = \log_e$
- $\log e = 1$ (because $e^1 = e$)
- $\log_{10} 1000 = 3$ (because $10^3 = 1000$)

Logistics
Splitting your data
Qualitative responses and Logistic Model

# A Brief Reminder About (Natural) Logarithms

Logarithm log is a class of functions.

- $\log_e z = x$ if $e^x = z$.
- We'll call $\log_e$ natural logarithm. And we'll assume $\log = \log_e$
- $\log e = 1$ (because $e^1 = e$)
- $\log_{10} 1000 = 3$ (because $10^3 = 1000$)

Some rules of logarithms

Logistics
Splitting your data
Qualitative responses and Logistic Model

# A Brief Reminder About (Natural) Logarithms

Logarithm log is a class of functions.

- $\log_e z = x$ if $e^x = z$.
- We'll call $\log_e$ natural logarithm. And we'll assume $\log = \log_e$
- $\log e = 1$ (because $e^1 = e$)
- $\log_{10} 1000 = 3$ (because $10^3 = 1000$)

Some rules of logarithms

- $\exp(\log(a)) = e^{\log(a)} = a$

Logistics
Splitting your data
Qualitative responses and Logistic Model

# A Brief Reminder About (Natural) Logarithms

Logarithm log is a class of functions.

- $\log_e z = x$ if $e^x = z$.
- We'll call $\log_e$ natural logarithm. And we'll assume $\log = \log_e$
- $\log e = 1$ (because $e^1 = e$)
- $\log_{10} 1000 = 3$ (because $10^3 = 1000$)

Some rules of logarithms

- $\exp(\log(a)) = e^{\log(a)} = a$
- $\log(a \times b) = \log(a) + \log(b)$ (!!!!!!)

Logistics
Splitting your data
Qualitative responses and Logistic Model

# A Brief Reminder About (Natural) Logarithms

Logarithm log is a class of functions.

- $\log_e z = x$ if $e^x = z$.
- We'll call $\log_e$ natural logarithm. And we'll assume $\log = \log_e$
- $\log e = 1$ (because $e^1 = e$)
- $\log_{10} 1000 = 3$ (because $10^3 = 1000$)

Some rules of logarithms

- $\exp(\log(a)) = e^{\log(a)} = a$
- $\log(a \times b) = \log(a) + \log(b)$ (!!!!!!)
- $\log(\frac{a}{b}) = \log(a) - \log(b)$

Logistics
Splitting your data
Qualitative responses and Logistic Model

# A Brief Reminder About (Natural) Logarithms

Logarithm log is a class of functions.

- $\log_e z = x$ if $e^x = z$.
- We'll call $\log_e$ natural logarithm. And we'll assume $\log = \log_e$
- $\log e = 1$ (because $e^1 = e$)
- $\log_{10} 1000 = 3$ (because $10^3 = 1000$)

Some rules of logarithms

- $\exp(\log(a)) = e^{\log(a)} = a$
- $\log(a \times b) = \log(a) + \log(b)$ (!!!!!!)
- $\log(\frac{a}{b}) = \log(a) - \log(b)$
- $\log(a^b) = b \log(a)$

Logistics
Splitting your data
Qualitative responses and Logistic Model

# A Brief Reminder About (Natural) Logarithms

Logarithm log is a class of functions.

- $\log_e z = x$ if $e^x = z$.
- We'll call $\log_e$ natural logarithm. And we'll assume $\log = \log_e$
- $\log e = 1$ (because $e^1 = e$)
- $\log_{10} 1000 = 3$ (because $10^3 = 1000$)

Some rules of logarithms

- $\exp(\log(a)) = e^{\log(a)} = a$
- $\log(a \times b) = \log(a) + \log(b)$ (!!!!!!)
- $\log(\frac{a}{b}) = \log(a) - \log(b)$
- $\log(a^b) = b\log(a)$
- $\log(1) = 0$

Logistics
Splitting your data
Qualitative responses and Logistic Model

# A Brief Reminder About (Natural) Logarithms

Logarithm log is a class of functions.

- $\log_e z = x$ if $e^x = z$.
- We'll call $\log_e$ natural logarithm. And we'll assume $\log = \log_e$
- $\log e = 1$ (because $e^1 = e$)
- $\log_{10} 1000 = 3$ (because $10^3 = 1000$)

Some rules of logarithms

- $\exp(\log(a)) = e^{\log(a)} = a$
- $\log(a \times b) = \log(a) + \log(b)$ (!!!!!!)
- $\log(\frac{a}{b}) = \log(a) - \log(b)$
- $\log(a^b) = b \log(a)$
- $\log(1) = 0$

Logistics
Splitting your data
Qualitative responses and Logistic Model

Call $p_i = \Pr(\text{Vote}_i = 1 | \boldsymbol{x}_i)$

Logistics
Splitting your data
Qualitative responses and Logistic Model

Call $p_i = \Pr(\text{Vote}_i = 1 | \boldsymbol{x}_i)$

$$\text{Vote}_i \quad \sim \quad \text{Bernoulli}(p_i)$$

Logistics
Splitting your data
Qualitative responses and Logistic Model

Call $p_i = \Pr(\text{Vote}_i = 1 | \boldsymbol{x}_i)$

$$\text{Vote}_i \sim \text{Bernoulli}(p_i)$$
$$p_i = f(\boldsymbol{\beta} \cdot \boldsymbol{x}_i)$$

Logistics
Splitting your data
Qualitative responses and Logistic Model

Call $p_i = \Pr(\text{Vote}_i = 1 | \mathbf{x}_i)$

$$\text{Vote}_i \sim \text{Bernoulli}(p_i)$$
$$p_i = f(\boldsymbol{\beta} \cdot \mathbf{x}_i)$$
$$\log\left(\frac{p_i}{1 - p_i}\right) = \boldsymbol{\beta} \cdot \mathbf{x}_i$$

Logistics
Splitting your data
Qualitative responses and Logistic Model

Call $p_i = \Pr(\text{Vote}_i = 1 | \mathbf{x}_i)$

$$
\begin{aligned}
\text{Vote}_i &\sim \text{Bernoulli}(p_i) \\
p_i &= f(\boldsymbol{\beta} \cdot \mathbf{x}_i) \\
\log\left(\frac{p_i}{1 - p_i}\right) &= \boldsymbol{\beta} \cdot \mathbf{x}_i \\
p_i &= \frac{\exp(\boldsymbol{\beta} \cdot \mathbf{x}_i)}{1 + \exp(\boldsymbol{\beta} \cdot \mathbf{x}_i)}
\end{aligned}
$$

Logistics
Splitting your data
Qualitative responses and Logistic Model

Call $p_i = \Pr(\text{Vote}_i = 1 | \boldsymbol{x}_i)$

$$
\begin{aligned}
\text{Vote}_i &\sim \text{Bernoulli}(p_i) \\
p_i &= f(\boldsymbol{\beta} \cdot \boldsymbol{x}_i) \\
\log\left(\frac{p_i}{1 - p_i}\right) &= \boldsymbol{\beta} \cdot \boldsymbol{x}_i \\
p_i &= \frac{\exp(\boldsymbol{\beta} \cdot \boldsymbol{x}_i)}{1 + \exp(\boldsymbol{\beta} \cdot \boldsymbol{x}_i)} \\
&= \frac{1}{1 + \exp(-\boldsymbol{\beta} \cdot \boldsymbol{x}_i)}
\end{aligned}
$$

Logistics
Splitting your data
Qualitative responses and Logistic Model

Call $p_i = \Pr(\text{Vote}_i = 1 | \boldsymbol{x}_i)$

$$
\begin{aligned}
\text{Vote}_i &\sim \text{Bernoulli}(p_i) \\
p_i &= f(\boldsymbol{\beta} \cdot \boldsymbol{x}_i) \\
\log\left(\frac{p_i}{1 - p_i}\right) &= \boldsymbol{\beta} \cdot \boldsymbol{x}_i \\
p_i &= \frac{\exp(\boldsymbol{\beta} \cdot \boldsymbol{x}_i)}{1 + \exp(\boldsymbol{\beta} \cdot \boldsymbol{x}_i)} \\
&= \frac{1}{1 + \exp(-\boldsymbol{\beta} \cdot \boldsymbol{x}_i)}
\end{aligned}
$$

Important functions:

Logistics
Splitting your data
Qualitative responses and Logistic Model

Call $p_i = \Pr(\text{Vote}_i = 1 | \boldsymbol{x}_i)$

$$
\begin{aligned}
\text{Vote}_i &\sim \text{Bernoulli}(p_i) \\
p_i &= f(\boldsymbol{\beta} \cdot \boldsymbol{x}_i) \\
\log\left(\frac{p_i}{1 - p_i}\right) &= \boldsymbol{\beta} \cdot \boldsymbol{x}_i \\
p_i &= \frac{\exp(\boldsymbol{\beta} \cdot \boldsymbol{x}_i)}{1 + \exp(\boldsymbol{\beta} \cdot \boldsymbol{x}_i)} \\
&= \frac{1}{1 + \exp(-\boldsymbol{\beta} \cdot \boldsymbol{x}_i)}
\end{aligned}
$$

Important functions:

$$
\text{odds}(p) = \frac{p}{1 - p}
$$

Logistics
Splitting your data
Qualitative responses and Logistic Model

Call $p_i = \Pr(\text{Vote}_i = 1 | \boldsymbol{x}_i)$

$$
\begin{aligned}
\text{Vote}_i &\sim \text{Bernoulli}(p_i) \\
p_i &= f(\boldsymbol{\beta} \cdot \boldsymbol{x}_i) \\
\log\left(\frac{p_i}{1 - p_i}\right) &= \boldsymbol{\beta} \cdot \boldsymbol{x}_i \\
p_i &= \frac{\exp(\boldsymbol{\beta} \cdot \boldsymbol{x}_i)}{1 + \exp(\boldsymbol{\beta} \cdot \boldsymbol{x}_i)} \\
&= \frac{1}{1 + \exp(-\boldsymbol{\beta} \cdot \boldsymbol{x}_i)}
\end{aligned}
$$

Important functions:

$$
\begin{aligned}
\text{odds}(p) &= \frac{p}{1 - p} \\
\text{log odds or logit}(p) &= \log\left(\frac{p}{1 - p}\right)
\end{aligned}
$$

Logistics
Splitting your data
Qualitative responses and Logistic Model

Call $p_i = \Pr(\text{Vote}_i = 1 | \boldsymbol{x}_i)$

$$
\begin{aligned}
\text{Vote}_i &\sim \text{Bernoulli}(p_i) \\
p_i &= f(\boldsymbol{\beta} \cdot \boldsymbol{x}_i) \\
\log\left(\frac{p_i}{1 - p_i}\right) &= \boldsymbol{\beta} \cdot \boldsymbol{x}_i \\
p_i &= \frac{\exp(\boldsymbol{\beta} \cdot \boldsymbol{x}_i)}{1 + \exp(\boldsymbol{\beta} \cdot \boldsymbol{x}_i)} \\
&= \frac{1}{1 + \exp(-\boldsymbol{\beta} \cdot \boldsymbol{x}_i)}
\end{aligned}
$$

Important functions:

$$
\begin{aligned}
\text{odds}(p) &= \frac{p}{1 - p} \\
\text{log odds or } \text{logit}(p) &= \log\left(\frac{p}{1 - p}\right) \\
\text{logistic function or } \text{logit}^{-1}(a) &= \frac{1}{1 + \exp(-a)}
\end{aligned}
$$

Logistics
Splitting your data
Qualitative responses and Logistic Model

# More on tutorial

Logistics
Splitting your data
Qualitative responses and Logistic Model

# R!

Logistics
Splitting your data
Qualitative responses and Logistic Model

# Some Context: Iraq Vote



- ▶ In 2002 President George Bush announced the Joint Resolution to Authorize the Use of United States Armed Forces Against Iraq.

- ▶ Congressional opposition.

Logistics
Splitting your data
Qualitative responses and Logistic Model

# NEXT

- Assessing model predictions for classification
  - Precision
  - Accuracy
  - Recall