

Introduction to Machine Learning for Social Scientists

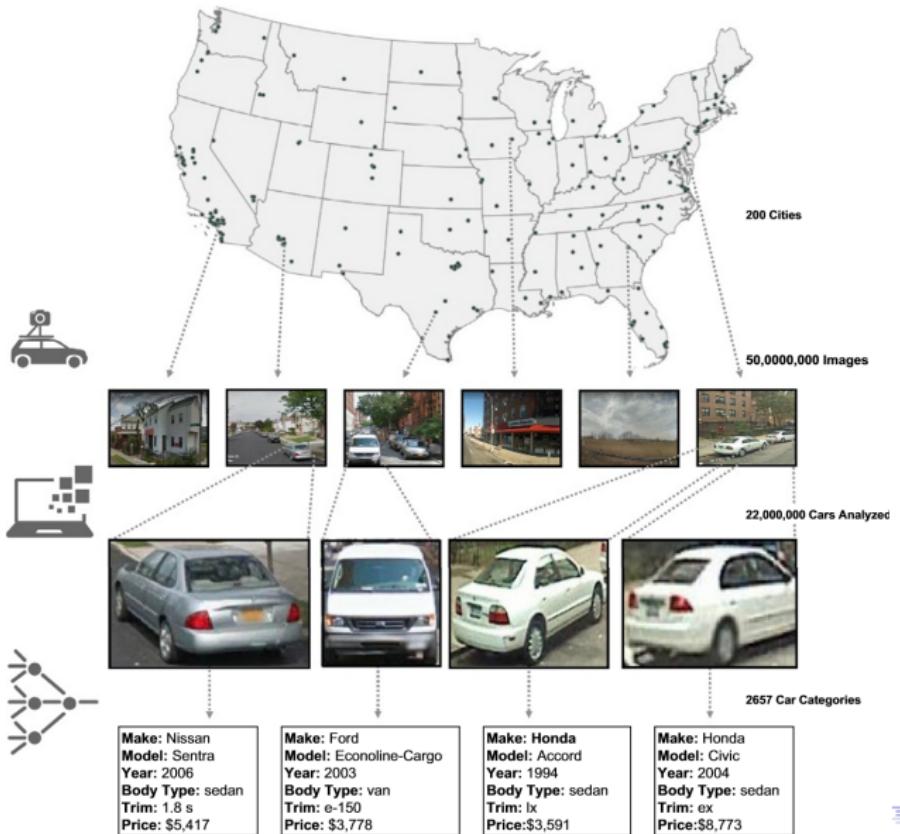
Class 7: More on performance and intro to re-sampling

Edgar Franco Vivanco

Stanford University
Department of Political Science

edgarf1@stanford.edu

Summer 2018



Homework 3: Due Friday July 20th at 11:59pm

Midterm: Monday July 23rd (1:30pm
to 2:50pm)

Submit your teams for the final project

By Wed 25th before class

One submission per team: spokesperson

Today

- ▶ Mid-term guidelines
- ▶ More on accuracy, precision and recall
- ▶ Intro to re-sampling

Midterm structure:

In-class, closed book.

Midterm structure:

In-class, closed book.

- ▶ General concepts (25%)
- ▶ Regression (25%)
- ▶ Classification (25%)
- ▶ R (25%)

Midterm structure:

In-class, closed book.

- ▶ General concepts (25%)
- ▶ Regression (25%)
- ▶ Classification (25%)
- ▶ R (25%)

No computers! No cellphone! Just a pen

General concepts:

- ▶ Inference vs. Prediction
- ▶ Supervised vs. Unsupervised
- ▶ Regression vs. Classification
- ▶ Training, test and validation sets
- ▶ Overfitting
- ▶ Cost functions, errors
- ▶ Models and variables (output, input, n)
- ▶ Ability to apply and interpret these concepts

Regression:

- ▶ Coefficients
- ▶ Prediction
- ▶ Limitations of linear model
- ▶ Errors
- ▶ SSE, RSS

Classification:

- ▶ Logistic model
- ▶ From probabilities to classes
- ▶ Confusion matrix
- ▶ Performance measures

R:

- ▶ Vector types (numeric, logical, string)
- ▶ Other basic objects (matrix, dataframe, list)
- ▶ Subset your data ([], \$, subset())
- ▶ Canned functions: predict, lm, glm, table, ifelse, set.seed, etc.
- ▶ Your own functions

A note on pseudocode:

```
#####
# Function to calculate pearson's correlation
# between two numeric vectors
# Syntax:
# my_cor(x,y):
# x: numeric vector of size n
# y: numeric vector of size n
#####
my_cor <- function(x,y){
  ## Calculating means
  mean_x <- mean(x)
  mean_y <- mean(y)
  ## Covariance
  cov_x_y <- sum((x-mean_x) * (y-mean_y))
  ## standard deviations
  sdx <- sum((x-mean_x)^2) * sum((y-mean_y)^2)
  ## Correlation: Cov/SD
  cor_x_y <- cov_x_y/sqrt(sdx)
  ### Return result
  return(cor_x_y)
}
```

A note on pseudocode:

```
#####
# Function to calculate pearson's correlation
# between two numeric vectors
# Syntax:
# my_cor(x,y):
# x: numeric vector of size n
# y: numeric vector of size n
#####
my_cor <- function(x,y){
  ## Calculating means
  mean_x <- mean(x)
  mean_y <- mean(y)
  ## Covariance
  cov_x_y <- sum((x-mean_x) * (y-mean_y))
  ## standard deviations
  sdx <- sum((x-mean_x)^2) * sum((y-mean_y)^2)
  ## Correlation: Cov/SD
  cor_x_y <- cov_x_y/sqrt(sdx)
  ### Return result
  return(cor_x_y)
}
```

my-corr \leftarrow function (X, Y){
1. Calculate mean of X
2. Calculate mean of Y
3. Using those means calculate Covariance
4. Using those means calculate Standard deviation
5. Divide Covariance over Standard deviation
6. Return result
}

A note on calculations:

$$SSE = \sum_{i=1}^N (Y_i - \hat{Y}_i)^2$$

$Y = c(2.6, 6.6, 5.5, 4, 3)$

$\hat{Y}_i = c(2, 5, 5.5, 4.1, 2.9)$

Where to study?

- ▶ Slides
- ▶ Code from class
- ▶ Section materials
- ▶ Homework solutions
- ▶ Book (ISL):
 - ▶ Chapter 2
 - ▶ Chapter 3 (3.1, 3.2, 3.6)
 - ▶ Chapter 4 (4.1, 4.2, 4.6)
 - ▶ Chapter 5 (5.1.1, 5.1.2)
- ▶ Extra tutorials

Questions?

Overview

Logistics

Midterm

More on Performance

Resampling

Confusion Matrix

To assess the quality of our data we compare our classifications with the real data or the "gold standard".

		Guess	
		Yes	No
Actual	Yes		
	No		

Confusion Matrix

Actual \ Guess		Yes	No
Actual	Yes	True positive	False Negative
No	False Positive	True Negative	

Accuracy

Accuracy is the percentage of observations classified correctly.

$$\text{Accuracy} = \frac{\text{TruePositive} + \text{TrueNegative}}{\text{TruePositive} + \text{TrueNegative} + \text{FalseNegative} + \text{FalsePositive}}$$

Precision

How many items classified as Yes are correctly classified?

$$Precision = \frac{TruePositive}{TruePositive + FalsePositive}$$

It is equal to 1 if all the guesses as Yes are actually Yes.

Recall

How many items **that are actually** as Yes are correctly classified?
In other words, is the number of correct results divided by the
number of results that should have been returned.

$$\text{Recall} = \frac{\text{TruePositive}}{\text{TruePositive} + \text{FalseNegative}}$$

It is equal to 1 if all the actual Yes are classified as Yes.

F-score

Harmonic mean of precision and recall:

$$F = \frac{2 * Precision * Recall}{Precision + Recall}$$

Example:

		Guess	
		Yes	No
Actual	Yes	10	15
	No	25	100

$$\text{Accuracy} = \frac{10+100}{(10+15+25+100)} = 0.73$$

Example:

Compared to a dumb classifier: All as negative

		Guess	
		Yes	No
Actual	Yes	0	25
	No	0	125

$$\text{Accuracy} = 0+125 / (0+25+0+125) = 0.83$$

This useless model has zero predictive power, yet better accuracy!

Other Example:

Let's try calculating accuracy for the following model that classified 100 tumors as malignant (the positive class) or benign (the negative class):

		Guess	Yes	No
		Actual		
		Yes	1	8
		No	1	90

- ▶ Accuracy = $1+90 / (1+90+8+1) = 0.91$
- ▶ Pretty good, right?

Other Example:

Let's try calculating accuracy for the following model that classified 100 tumors as malignant (the positive class) or benign (the negative class):

		Guess	
		Yes	No
Actual	Yes	1	8
	No	1	90

- ▶ Accuracy = $1+90 / (1+90+8+1) = 0.91$
- ▶ Pretty good, right?
- ▶ Of the 100 tumor examples, 91 are benign (90 TNs and 1 FP) and 9 are malignant (1 TP and 8 FNs).

Other Example:

Let's try calculating accuracy for the following model that classified 100 tumors as malignant (the positive class) or benign (the negative class):

		Guess	Yes	No
		Actual		
		Yes	1	8
		No	1	90

- ▶ Accuracy = $1+90 / (1+90+8+1) = 0.91$
- ▶ Pretty good, right?
- ▶ Of the 100 tumor examples, 91 are benign (90 TNs and 1 FP) and 9 are malignant (1 TP and 8 FNs).
- ▶ Of the 9 malignant tumors, the model only correctly identifies 1 as malignant— a terrible outcome, **as 8 out of 9 malignancies go undiagnosed!**

Precision and Recall:

Let's try calculating precision and recall for the following model that classified 100 tumors as malignant (the positive class) or benign (the negative class):

Actual	Guess	
	Yes	No
Yes	1	8
No	1	90

Precision and Recall:

Let's try calculating precision and recall for the following model that classified 100 tumors as malignant (the positive class) or benign (the negative class):

		Guess	
		Yes	No
Actual	Yes	1	8
	No	1	90

- ▶ Precision = $1 / (1+1) = 0.5$
- ▶ Our model has a precision of 0.5 – in other words, when it predicts a tumor is malignant, it is correct 50% of the time.

Precision and Recall:

Let's try calculating precision and recall for the following model that classified 100 tumors as malignant (the positive class) or benign (the negative class):

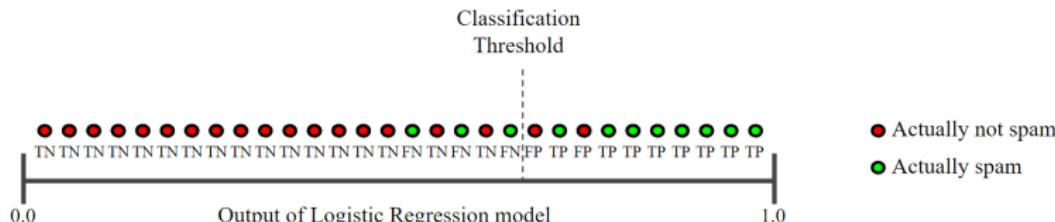
Actual \ Guess	Yes	No
Yes	1	8
No	1	90

- ▶ Precision = $1 / (1+1) = 0.5$
- ▶ Our model has a precision of 0.5 – in other words, when it predicts a tumor is malignant, it is correct 50% of the time.
- ▶ Recall = $1 / (1+8) = 0.11$
- ▶ Our model has a recall of 0.11 – in other words, it correctly identifies 11% of all malignant tumors.

Precision and Recall:

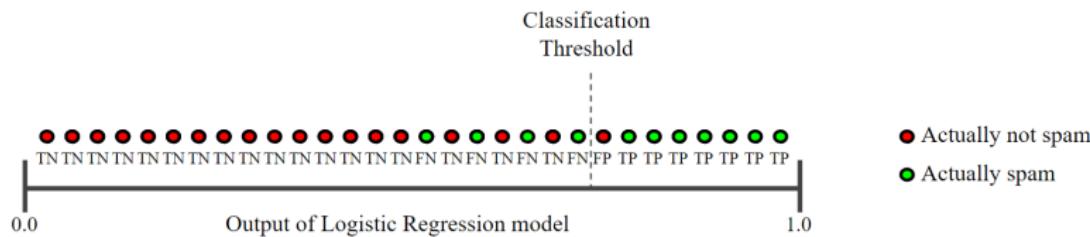
- ▶ Precision and Recall are often in tension
- ▶ Example: Classifying Spam
- ▶ Precision: $8 / (8+2) = .8$
- ▶ Recall: $8 / (8+3) = 0.72$

Guess		Yes	No
Actual	Yes	8	3
No	2	17	



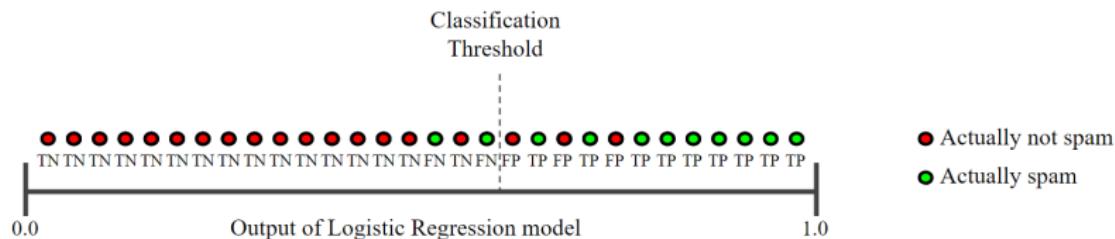
The effect of increasing the classification threshold.:

- ▶ The number of false positives decreases, but false negatives increase.
- ▶ Precision: $7 / (7+1) = 0.88$
- ▶ Recall: $7 / (7+4) = 0.64$
- ▶ As a result, precision increases, while recall decreases.



The effect of decreasing the classification threshold.:

- ▶ False positives increase, and false negatives decrease.
- ▶ Precision: $9 / (9+3) = 0.75$
- ▶ Recall: $9 / (9+2) = 0.82$
- ▶ As a result, this time, precision decreases and recall increases.



Some lessons

- ▶ Accuracy alone doesn't tell the full story, in particular when you're working with a class-imbalanced data set.

Some lessons

- ▶ Accuracy alone doesn't tell the full story, in particular when you're working with a class-imbalanced data set.
- ▶ Usually, to fully evaluate the effectiveness of a model, you must examine both precision and recall.

Some lessons

- ▶ Accuracy alone doesn't tell the full story, in particular when you're working with a class-imbalanced data set.
- ▶ Usually, to fully evaluate the effectiveness of a model, you must examine both precision and recall.
- ▶ There's always a trade-off between precision and recall.

Some lessons

- ▶ Accuracy alone doesn't tell the full story, in particular when you're working with a class-imbalanced data set.
- ▶ Usually, to fully evaluate the effectiveness of a model, you must examine both precision and recall.
- ▶ There's always a trade-off between precision and recall.
- ▶ At the end, there is not a best way to evaluate all models.

Some lessons

- ▶ Accuracy alone doesn't tell the full story, in particular when you're working with a class-imbalanced data set.
- ▶ Usually, to fully evaluate the effectiveness of a model, you must examine both precision and recall.
- ▶ There's always a trade-off between precision and recall.
- ▶ At the end, there is not a best way to evaluate all models.
- ▶ Different metrics give us different (and valuable) insights into how a classification model performs.

Overview

Logistics

Midterm

More on Performance

Resampling

Evaluating fit:

- ▶ In sample: dependent variable in training data

Evaluating fit:

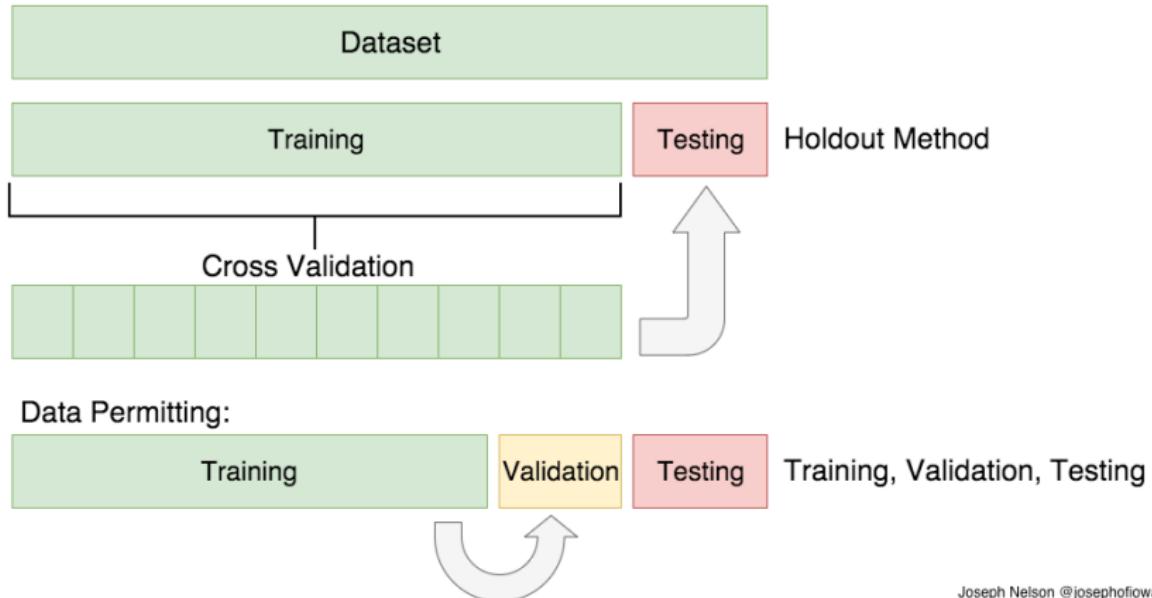
- ▶ In sample: dependent variable in training data
- ▶ But in general, we do not really care how well the method works on the training data. Rather, we are interested in the accuracy of the predictions that we obtain when we apply our method to previously unseen test data.

Evaluating fit:

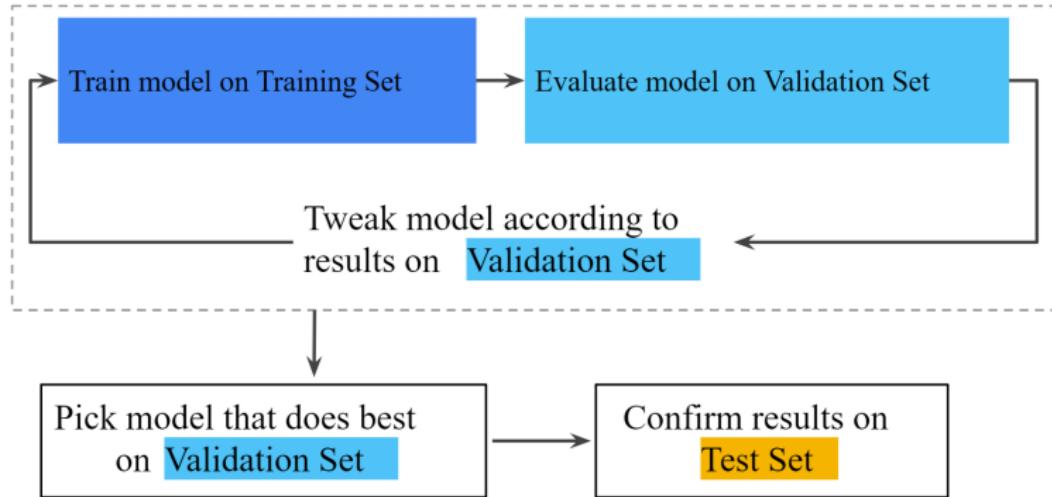
- ▶ In sample: dependent variable in training data
- ▶ But in general, we do not really care how well the method works on the training data. Rather, we are interested in the accuracy of the predictions that we obtain when we apply our method to previously unseen test data.
- ▶ Out of sample: **held out** data, test data

Evaluating fit:

- ▶ In sample: dependent variable in training data
- ▶ But in general, we do not really care how well the method works on the training data. Rather, we are interested in the accuracy of the predictions that we obtain when we apply our method to previously unseen test data.
- ▶ Out of sample: **held out** data, test data
- ▶ Best practice: evaluate t with **gold standard** data



Joseph Nelson @josephofiowa



Cross-Validation: Some Intuition

Recall Optimal division of data:

Cross-Validation: Some Intuition

Recall Optimal division of data:

- Train: build model

Cross-Validation: Some Intuition

Recall Optimal division of data:

- Train: build model
- Validation: assess model

Cross-Validation: Some Intuition

Recall Optimal division of data:

- Train: build model
- Validation: assess model
- Test: classify remaining documents

Cross-Validation: Some Intuition

Recall Optimal division of data:

- Train: build model
- Validation: assess model
- Test: classify remaining documents

K-fold Cross-validation idea: create many training and test sets.

Cross-Validation: Some Intuition

Recall Optimal division of data:

- Train: build model
- Validation: assess model
- Test: classify remaining documents

K-fold Cross-validation idea: create many training and test sets.

- Idea: use observations both in training and test sets

Cross-Validation: Some Intuition

Recall Optimal division of data:

- Train: build model
- Validation: assess model
- Test: classify remaining documents

K-fold Cross-validation idea: create many training and test sets.

- Idea: use observations both in training and test sets
- Each step: use held out data to evaluate performance

Cross-Validation: Some Intuition

Recall Optimal division of data:

- Train: build model
- Validation: assess model
- Test: classify remaining documents

K-fold Cross-validation idea: create many training and test sets.

- Idea: use observations both in training and test sets
- Each step: use held out data to evaluate performance
- **Avoid overfitting** and have context specific penalty

Cross-Validation: Some Intuition

Recall Optimal division of data:

- Train: build model
- Validation: assess model
- Test: classify remaining documents

K-fold Cross-validation idea: create many training and test sets.

- Idea: use observations both in training and test sets
- Each step: use held out data to evaluate performance
- **Avoid overfitting** and have context specific penalty

Steps

1. Randomly divide the data into training and validation

Steps

1. Randomly divide the data into training and validation
2. Fit the model on the training set

Steps

1. Randomly divide the data into training and validation
2. Fit the model on the training set
3. The fitted model is used to predict the observations in the validation set

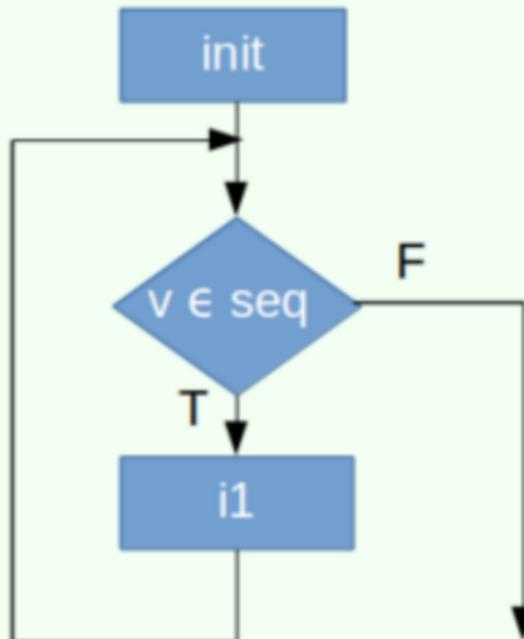
Steps

1. Randomly divide the data into training and validation
2. Fit the model on the training set
3. The fitted model is used to predict the observations in the validation set
4. The resulting validation set error rate typically assessed using MSE in the case of a quantitative response provides an estimate of the test error rate.

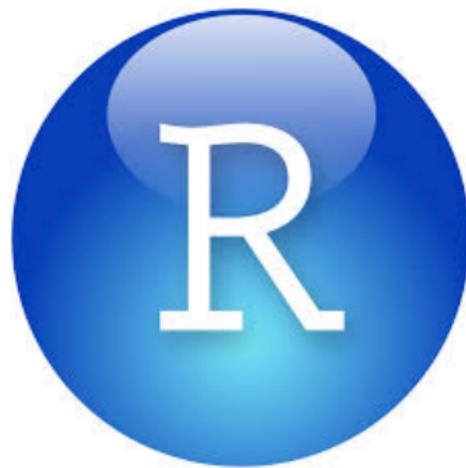
Leave-One-Out Cross-Validation

- ▶ Like the validation set approach, LOOCV involves splitting the set of observations into two parts.
- ▶ However, instead of creating two subsets of comparable size, a single observation (x_1, y_1) is used for the validation set, and the remaining observations $(x_2, y_2), \dots, (x_n, y_n)$ make up the training set.
- ▶ The statistical learning method is fit on the $n - 1$ training observations, and a prediction \hat{y}_1 is made for the excluded observation, calculating the error $(y_1 - \hat{y}_1)^2$
- ▶ We repeat the process n times
- ▶ We estimate a mean-cross validation error.

For loop



R!



Remember:

- ▶ Midterm (Study!)
- ▶ Homework
- ▶ Submit groups