# CMT307: Applied ML Session 3

Experimental design

*Jose Camacho Collados*
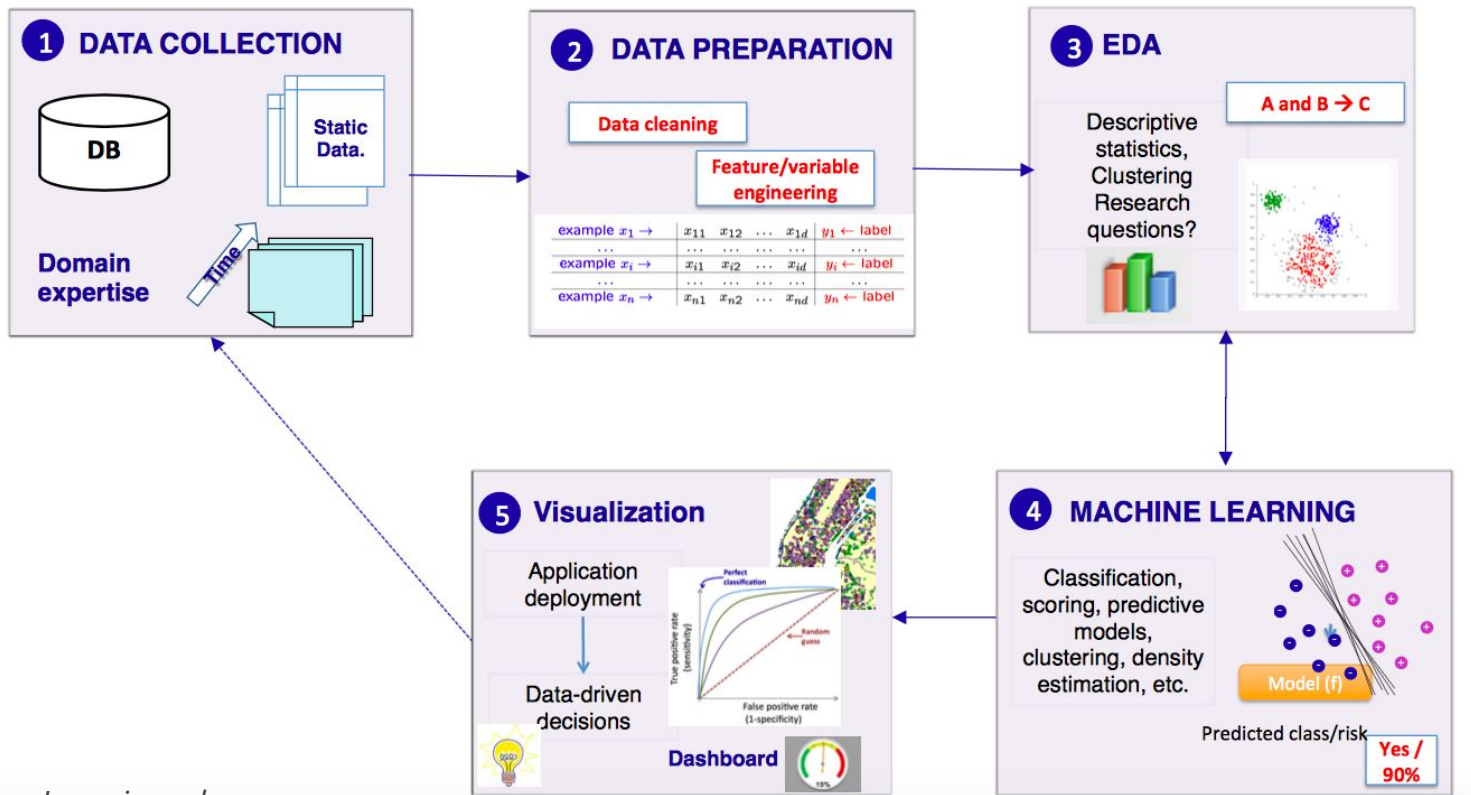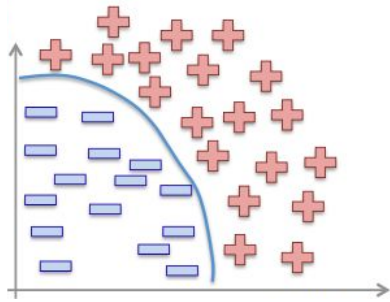
# Reminders (last session)

- Machine learning pipeline + Feature engineering/selection/extraction.

- Python notebook on basic **feature engineering/selection** with *sklearn* (Solutions to the exercises now available in learning central).
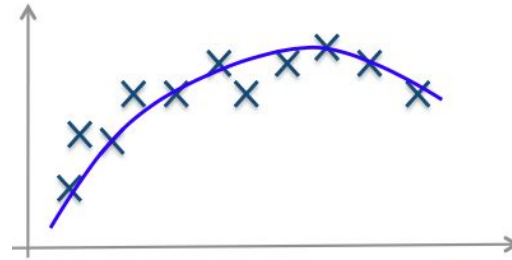
# Machine Learning pipeline (last sesssion)

3

# Classification vs. Regression (last session)

**Classification (binary)**
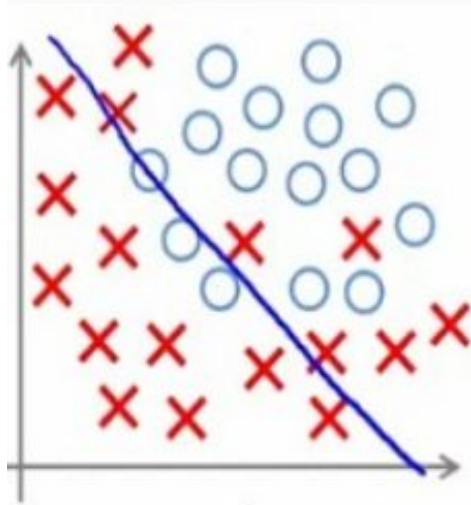
**Regression**

Categories

Continuous values
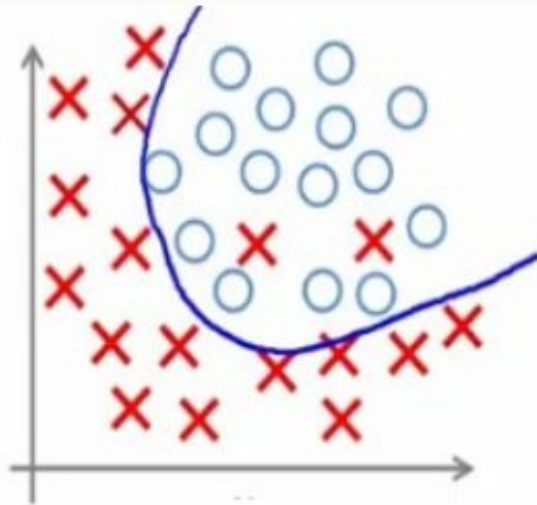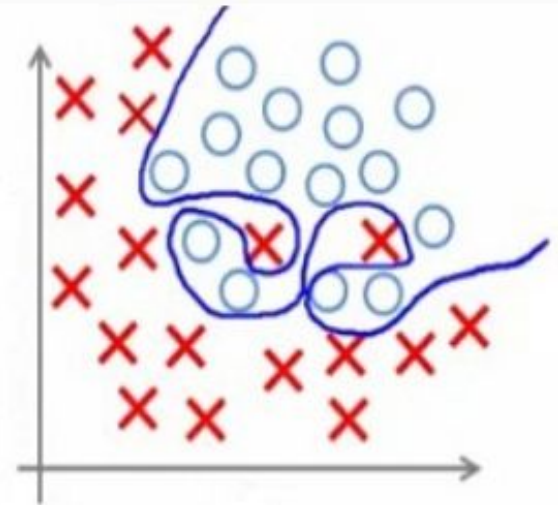
# Overfitting vs. underfitting (last session)



**Under-fitting**   **Appropriate-fitting**   **Over-fitting**

# Outline

➢ Housekeeping/opportunities

➢ Experimental design

➢ Evaluation measures

➢ Coursework

➢ Hands on!

# Panopto: Video-recorded lectures

Currently not working properly (IT issue)

It should be fixed for this week (hopefully!)

Two first sessions may be lost - apologies!

# Data and Knowledge Engineer Seminars

**Weekly seminars** (on Mondays from 1pm to 2pm) about machine learning, knowledge representation and data mining.

**Location:** Room WX/3.07 (West extension, Queen's building)

➢ **Last:** *Interface between AI and Social Sciences and why it matters*

➢ **Next (28 Oct):** *Sentence embeddings (NLP) + Classifying memory reactivation during sleep using machine learning*

*Link:*
https://www.cardiff.ac.uk/computer-science/events/data-and-knowledge-engineering-seminars

# Opportunities in Cardiff: China Scholarship Council

**Tuition fee PhD scholarships** in the fields of science and technology

How to apply: CV+Personal Statement

**First-Come-First-Served!** Apply early = More chances

**Application deadline:** 15 January 2020. **Start date**: October 2020

*[Only for Chinese students]*

*Link:*
https://www.cardiff.ac.uk/study/international/funding-and-fees/international-scholarships/phd-research-scholarships

# Experimental design

# Train, validation and test



Datasets are usually split into three parts:

➢ **Training set:** Set of instances (examples) used to train a model.

➢ **Development (validation) set:** Instances that are used to train the parameters of your model -> Useful to avoid overfitting!

➢ **Test set:** Split where the model is evaluated (only used once!).

# Train, validation and test



How to split the dataset?

**No general rule.** Generally training data larger, but test data should be large enough for the evaluation to be meaningful (and development for the model to generalize well).

Some usual train/dev/test splits: 80%/10%/10% or 60%/20%/20%
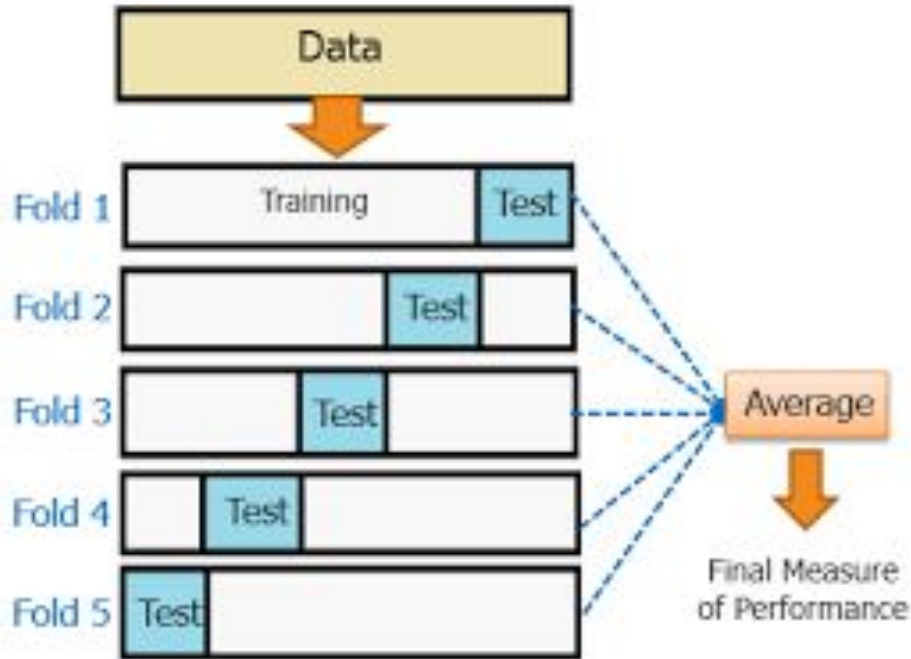
# Cross-validation

Automatic creation of train/test splits (dev if needed) using $k$ folds.

Useful when the dataset is not large, or with a very skewed distribution.

**Important:** No preprocessing or feature selection can be done in the initial full dataset! All steps must be performed within the folds.

# Example: 5-fold cross-validation

# Evaluation measures

# Evaluation measures (binary classification)

Binary classifiers are evaluated by comparing their predictions to a reference label set, also referred to as ground truth or **gold standard**.

The results of a classifier and an understanding of its behaviour can be easily condensed in a **confusion matrix** (or contingency table).

# Evaluation measures (binary classification)
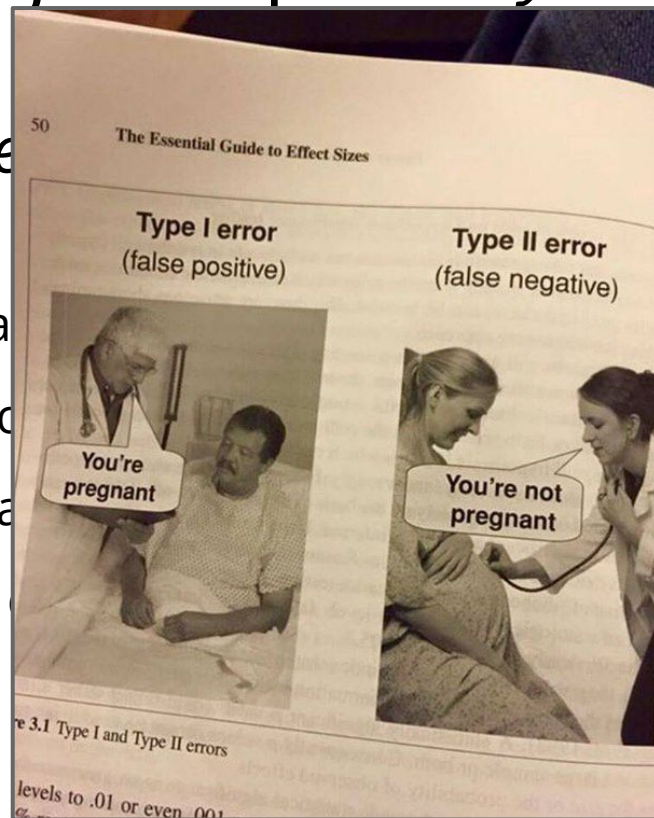
*Example: Predicting whether a person has diabetes or not*

➢ **True Positive (TP):** Our prediction came out positive, and it was right

➢ **True Negative (TN):** Our prediction came out negative, and it was right.

➢ **False Positive (FP):** Our prediction came out positive, but it was wrong.

➢ **False Negative (FN):** Our prediction came out negative, but it was wrong.

# Evaluation measures (binary classification)

*Example: Predicting whether a pe...* ...t

- ➢ **True Positive (TP):** Our prediction ca... ...ght

- ➢ **True Negative (TN):** Our prediction ... ...right.

- ➢ **False Positive (FP):** Our prediction ca... ...rong.

- ➢ **False Negative (FN):** Our prediction ... ...wrong.



18

# Confusion matrix

| n=165 | Predicted: NO | Predicted: YES | |
|---|---|---|---|
| Actual: NO | TN = 50 | FP = 10 | 60 |
| Actual: YES | FN = 5 | TP = 100 | 105 |
| | 55 | 110 | |

**TN:** True negative

**FN:** False negative

**FP:** False positive

**TP:** True positive

# Evaluation measures (binary classification)

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

| | | |
|---|---|---|
| **Accuracy** | (TP + TN) / (TP + TN + FP + FN) | The percentage of predictions that are correct |
| **Precision** | TP / (TP + FP) | The percentage of positive predictions that are correct |
| **Sensitivity (Recall)** | TP / (TP + FN) | The percentage of positive cases that were predicted as positive |
| **Specificity** | TN / (TN + FP) | The percentage of negative cases that were predicted as negative |

Slide credit: *Machine Learning edx*

# Evaluation measures: F1-Score

While accuracy gives us the percentage of correct answers, sometimes this is not enough (e.g. a dataset heavily biased towards one class).

Another important measure if **F1-Score** (or F-Measure), which is the harmonic mean of precision and recall.

$$F_1 = 2 * \frac{precision * recall}{precision + recall}$$

# Evaluation measures (multi-class classification)

➢ **Micro-average:** compute the metric independently for each class and then take the average (i.e. treat all the classes equally)

➢ **Macro-average:** aggregate the contributions of all classes to compute the average metric (i.e. more frequent classes have more weight).

Micro-average useful for general **overview**, and macro-average useful if we want our classifier to **work well in all classes**.

We can then compute **micro/macro precision, recall or F1-score**.

**More information:** http://rushdishams.blogspot.com/2011/08/micro-and-macro-average-of-precision.html

# Statistical significance tests

Many times we want to compare the performance between ML models. However, depending on the size of the dataset, the difference between two different models may be due to some statistical artifact.

To know how meaningful these differences are, we run **statistical hypothesis tests**.

Depending on the the case and evaluation metric, we may use different tests (**t-test**, **chi-squared, McNemar**, etc.).

**More information:**
https://machinelearningmastery.com/statistical-significance-tests-for-comparing-machine-learning-algorithms

# Evaluation measures (regression)

| | |
|---|---|
| Mean squared error | $\text{MSE} = \dfrac{1}{n} \sum_{t=1}^{n} e_t^2$ |
| Root mean squared error | $\text{RMSE} = \sqrt{\dfrac{1}{n} \sum_{t=1}^{n} e_t^2}$ |
| Mean absolute error | $\text{MAE} = \dfrac{1}{n} \sum_{t=1}^{n} |e_t|$ |
| Mean absolute percentage error | $\text{MAPE} = \dfrac{100\%}{n} \sum_{t=1}^{n} \left| \dfrac{e_t}{y_t} \right|$ |

Number of instances=$n$

"*e*" refers to the error in one particular instance (generally, the difference between the prediction and the gold label)

Image credit: *Vimarsh Karbhari*

24

# Evaluation measures (regression)

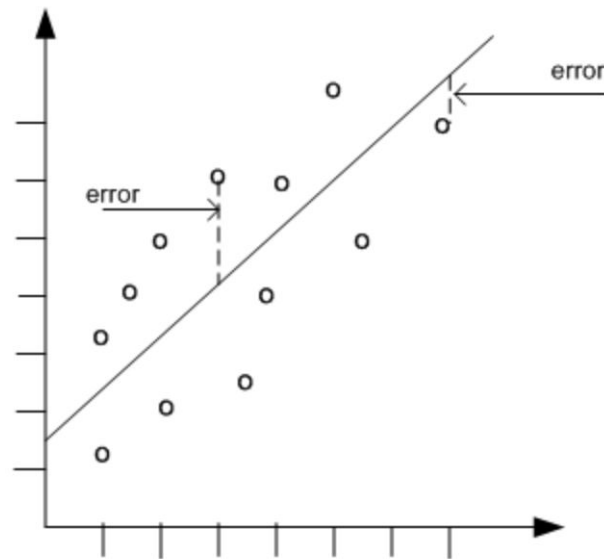| | |
|---|---|
| Mean squared error | $\mathrm{MSE} = \dfrac{1}{n} \sum_{t=1}^{n} e_t^2$ |
| Root mean squared error | $\mathrm{RMSE} = \sqrt{\dfrac{1}{n} \sum_{t=1}^{n} e_t^2}$ |
| Mean absolute error | $\mathrm{MAE} = \dfrac{1}{n} \sum_{t=1}^{n} |e_t|$ |
| Mean absolute percentage error | $\mathrm{MAPE} = \dfrac{100\%}{n} \sum_{t=1}^{n} \left| \dfrac{e_t}{y_t} \right|$ |



Image credit: *Vimarsh Karbhari*

# Coursework

# Coursework (1st semester)

➢ **Available in Learning Central:** ~Monday, October 28th

➢ **Submission date and time:** Tuesday, January 14th at 9:30am

Most of the exercises are highly practical (Python, sklearn, etc.)

Extra credit possible (a bit more challenging).

# School's private Stack Overflow (reminder)

https://stackoverflow.com/c/comsc



Post your questions related to the course here!

Add the tags **cmt307** and **machine-learning** to your question.

# Next sessions

**Next session:** 7 November

Following four sessions (Session 4-7) with Dr. Yuhua Li.

In these sessions you will get to understand and delve into the actual machine learning algorithms (SVM, random forests, etc.). A bit more mathematical.

**Session 8:** Guest lecture + ethics and bias + recap

See you in Session 8!

# Hands on!



Python notebook with exercises
about **experimental design** available at Learning Central