

Applied Machine Learning (AML)

Class Starting at 4:10pm

Oisin Mac Aodha • Siddharth N.

Applied Machine Learning

Week 7: Evaluation and Model Selection

This slides will be made available on the project website after the class. This session will be recorded.

Overview

- 1) Outline your tasks this for week
- 2) Discussion of Week 6's topics

Coursework Progress Report

- Progress reports due by 5pm, 29 Oct (Wed) (not assessed)
- Submit through form linked on website
 - Only one team member needs to submit
 - Use 2-digit 0-padded naming convention [e.g. 06.pdf, 23.pdf]
- (Optional) Feedback session 1-4pm, 31 Oct (Fri); AT 5.04
 - Primarily to check you are not miscalibrated
 - Will communicate schedule once we know signup

Coursework Final Submission

- Final submission due 12pm, 20 Nov (Thu)
 - Only report submission required here
- Supplementary materials deadline in the following week
 - Source for project (Readme + minimal structure)
 - Source for report
- Further details for both announced later

Coursework: Communication & Contribution

- Important to communicate with teammates do not neglect / delay
- Regular contact ensures everyone is on the same page
- Individual grade based on final report grade **
 - Statement of contribution required (not counted for page limit)
 - Disparities will be taken into account for individual grades
- Also remember Guidance on use of Generative Al
 - Must explicitly acknowledge use and describe how it was used

Week 7: Your tasks for this week

- 1) Complete Lab 3
- 2) Watch videos for week 7
 - Clustering and Non-Linear Dimensionality Reduction
- 3) Ask questions on Piazza if stuck
- 4) Continue working on the coursework
- 5) Start Tutorial 3 which takes places next week link in week 8

Evaluation

Which evaluation metrics are commonly used for evaluating the performance of a classification model when there is a class imbalance problem?

- 1. Accuracy and Error
- 2. Precision and Error
- 3. Precision and Recall
- 4. Error and Recall



Match the term to the formulation

- 1. Precision
- 2. Recall
- 3. False Positive Rate
- 4. True Positive Rate
- 5. False Negative Rate

$$A. \quad \frac{\text{TP}}{\text{TP} + \text{FP}}$$

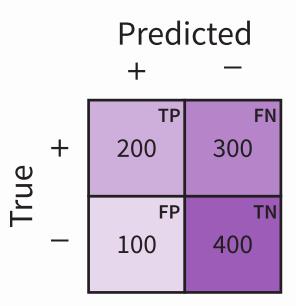
$$B. \quad \frac{\text{FP}}{\text{FP} + \text{TN}}$$

$$C. \quad \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$D. \quad \frac{\text{FN}}{\text{TP} + \text{FN}}$$

Error Measures

- False Positive Rate (FPR) = $\frac{FP}{FP + TN}$
 - (False Alarm) % of '-' misclassified as '+'
- False Negative Rate (FNR) = $\frac{FN}{TP + FN}$
 - (Miss) % of '+' misclassified as '-'
- Recall / True Positive Rate (TPR) = $\frac{TP}{TP + FN}$
 - (1 Miss) % of '+' correctly predicted
- Precision / Positive Predictive Rate (PPR) = $\frac{TP}{TP + FP}$
 - % of '+' out of all positive predictions

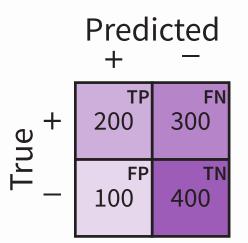


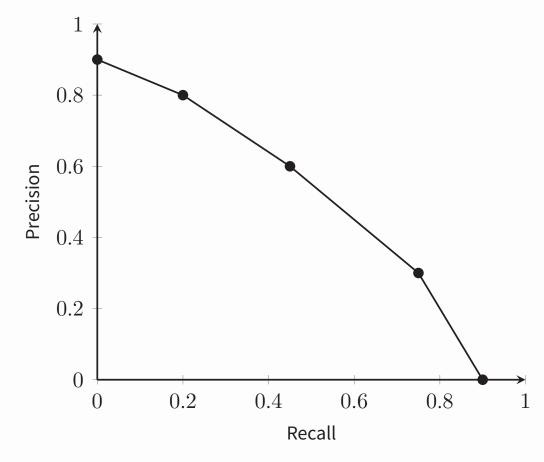


Precision-Recall Curve

$$Precision = \frac{TP}{TP + FP} \qquad Recall = \frac{TP}{TP + FN}$$

| email | label | $p(y \mathbf{x})$ |
|-------------------------|-------|-------------------|
| "send us your password" | + | 0.92 |
| "send us review" | _ | 0.80 |
| "review your account" | _ | 0.72 |
| "review us" | + | 0.65 |
| "send your password" | + | 0.61 |
| "send us your account" | + | 0.43 |
| • • | | |







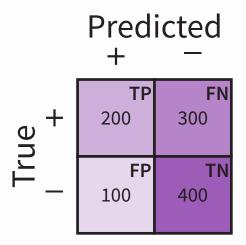
Receiver Operating Characteristic (ROC)

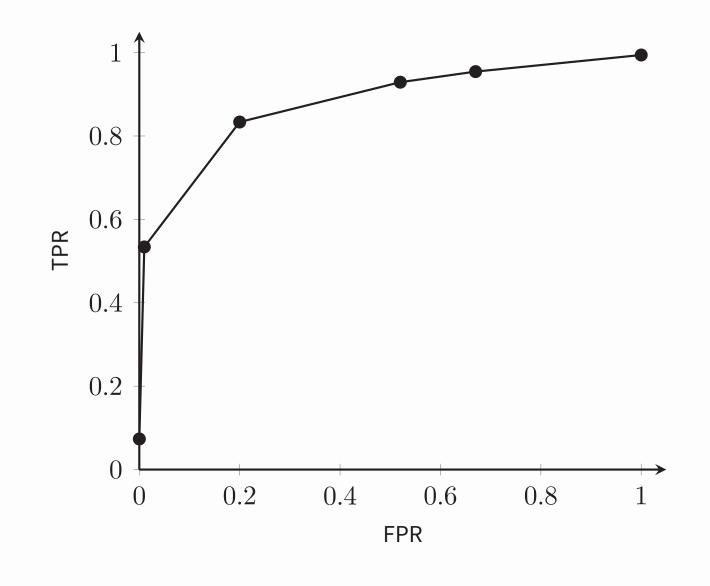
True Postive Rate (TPR) / Recall =
$$\frac{TP}{TP + FN}$$
 False Po

False Positive Rate (FPR) =
$$\frac{FP}{FP + TN}$$

AUC

- area under ROC curve
- larger area ⇒ better model



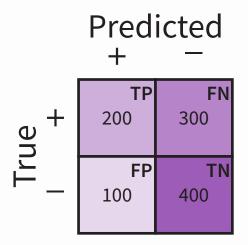


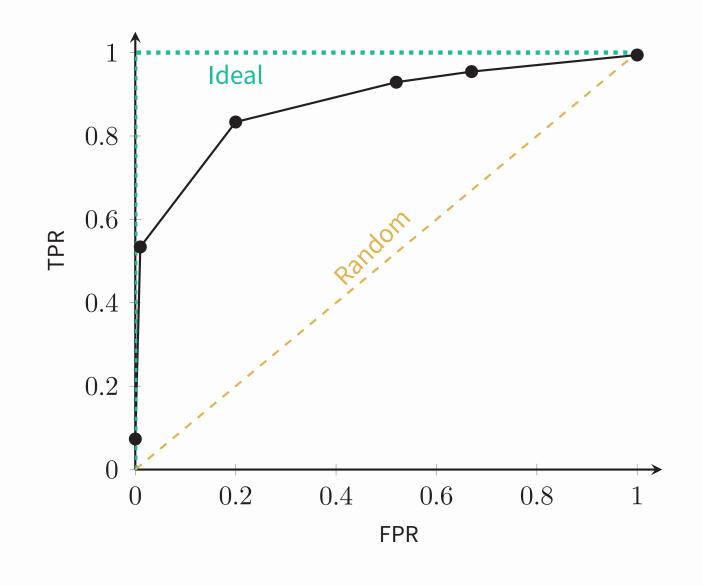
Receiver Operating Characteristic (ROC)

True Postive Rate (TPR) / Recall =
$$\frac{\mathrm{TP}}{\mathrm{TP} + \mathrm{FN}}$$
 False Positive Rate (FPR) = $\frac{\mathrm{FP}}{\mathrm{FP} + \mathrm{TN}}$

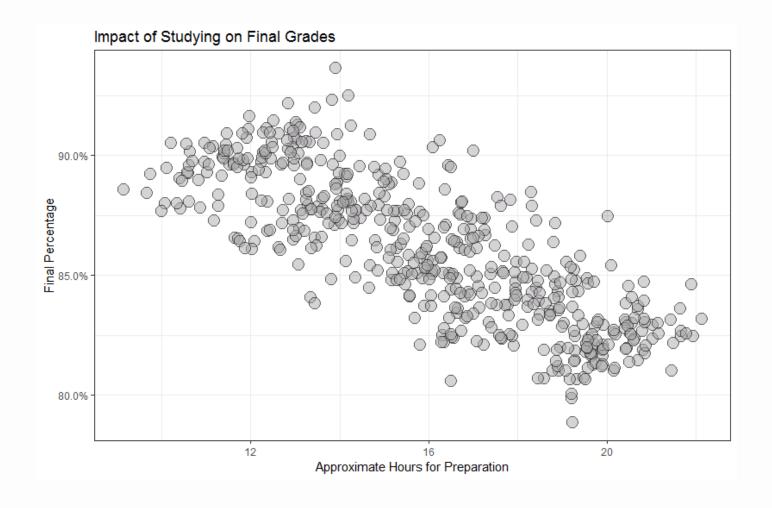
AUC

- area under ROC curve
- larger area ⇒ better model



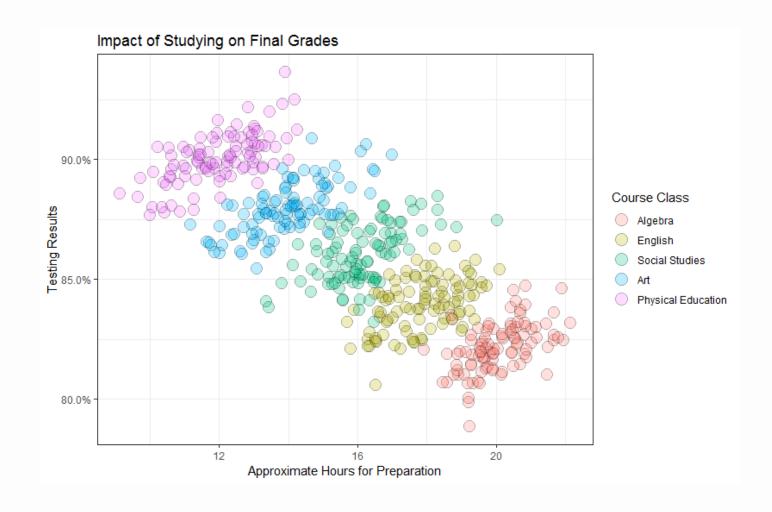


Simpson's Paradox





Simpson's Paradox





Model Selection

Preliminaries

Population vs. Sample statistics

Population: All the elements from a set

E.g. All leave-1-out splits of the dataset

Sample: Observations drawn from population

E.g. Some N splits of the dataset

If sample set is x_1, \ldots, x_N

$$\overline{x} = \frac{1}{N} \sum_{i=1}^{N} x_i$$

$$s^{2} = \frac{1}{N-1} \sum_{i=1}^{N} (x_{i} - \overline{x})^{2}$$

*Bessel's correction

Preliminaries

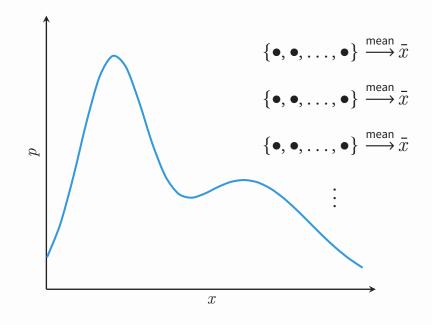
Central Limit Theorem (CLT)

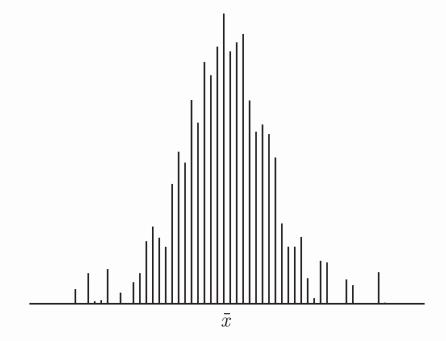
For a set of samples x_1, \ldots, x_N, \ldots from a population with expected mean μ and finite variance σ^2

$$z = \frac{\bar{x} - \mu}{\sigma / \sqrt{N}} \sim \mathcal{N}(0, 1) \quad \text{as } N \to \infty$$

Assume

- population μ known
- population σ^2 known







Preliminaries

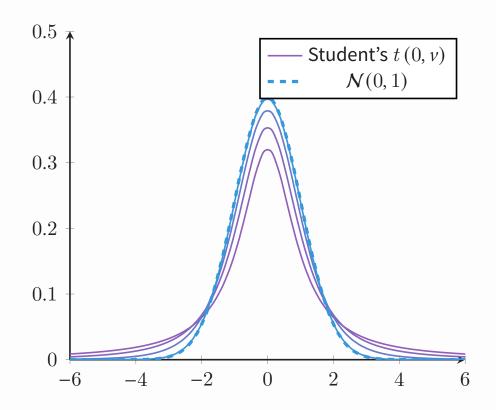
Student's-t distribution

- CLT: (weak) convergence to $\mathcal{N}(0,1)$ as $N \to \infty$
- for smaller *N*, not Gaussian!

Assume

- population μ known
- population σ^2 unknown
- estimate sample variance $s^2 = \frac{1}{N-1} \sum_{i=1}^{N} (x_i \overline{x}_N)^2$

$$t = \frac{\overline{x} - \mu}{s/\sqrt{N}}, \quad \nu = N - 1$$

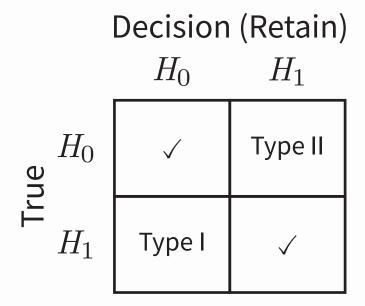


$$f(t, \nu) = \frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\nu\pi} \Gamma(\frac{\nu}{2})} \left(1 + \frac{t^2}{\nu}\right)^{-(\nu+1)/2}$$



Hypothesis Testing

- Formally examine two opposing conjectures (hypothesis): H_0 and H_1
- Mutually exclusive and exhaustive: $H_0 = \text{True} \implies H_1 = \text{False}$
- Analyse data to determine which is True and which is False



Null Hypothesis: H_0

- States the assumption to be tested
- Begin with assumption that $H_0 = \text{True}$
- Always evaluates (partial) equality $(=, \leq, \geq)$

Alternative Hypothesis: H_1

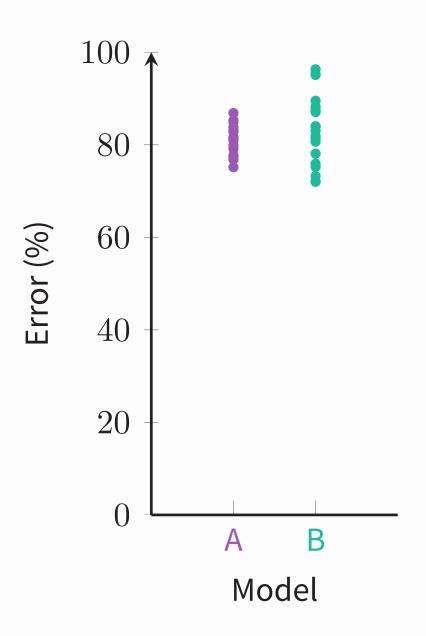
- States the assumption believed to be True
- Evaluate if evidence supports assumption
- Always evaluates (strict) inequality (≠, >, <)

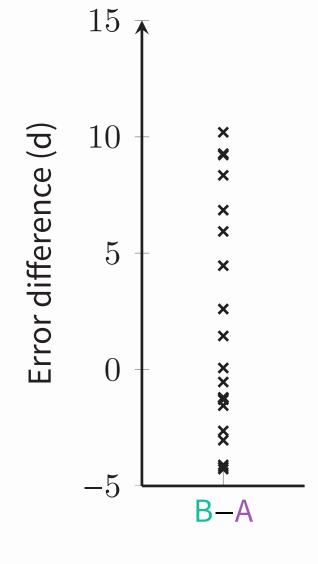


Example: Hypothesis Testing for Models

Generating Variation

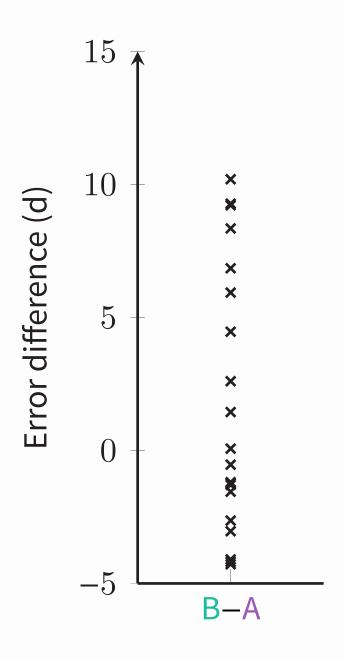
| Data Split | Α | В | d |
|---|-----------------------|------------|---------------------------|
| $\left\{\mathcal{D}_{train}^1, \mathcal{D}_{test}^1 \right\}$ | ℓ_1^A | ℓ_1^B | $\ell_1^B - \ell_1^A$ |
| $\left\{\mathcal{D}^2_{train}, \mathcal{D}^2_{test} \right\}$ | ℓ_2^A | ℓ_2^B | $\ell_2^B - \ell_2^A$ |
| • • | • | • | • • |
| $\left\{\mathcal{D}_{train}^{N}, \mathcal{D}_{test}^{N} ight\}$ | $oldsymbol{\ell}_N^A$ | ℓ_N^B | $\ell_N^{B} - \ell_N^{A}$ |







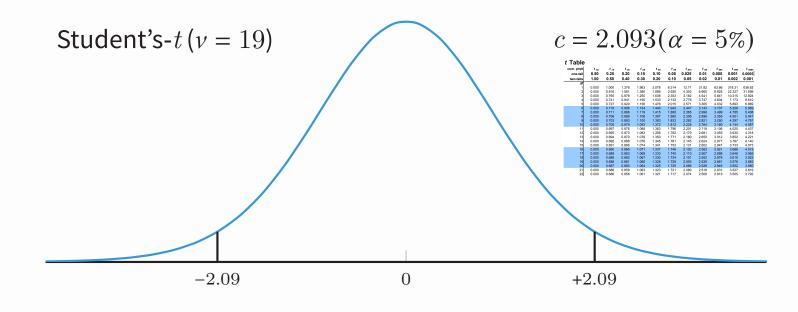
Example: Hypothesis Testing for Models



Hypotheses

 $H_0: \mu^d = 0$ $\alpha = 5\%$ (significance)

 $H_1: \mu^d \neq 0 \qquad \qquad N = 20$





Hypothesis Testing: Caveats

- Rejecting H_0 does not imply 100% sure H_0 is False
- ullet Failing to reject H_0 does not imply H_0 is True
- Confidence level ($\alpha = 0.05$) is from convention; not always best
- Statistical significance does not imply practical relevance
 - Rejecting $H_0: \mu^d = 0$ only tells us that $\mu^d \neq 0$ but not how big or important the difference is
 - Remedy: Report confidence interval (CI)

$$\bar{d} \pm c|_{\alpha} \cdot \frac{s}{\sqrt{N}}$$

which, for our example would be

$$2.53 \pm 2.093 \cdot \frac{5.27}{\sqrt{20}} \qquad (\alpha = 0.05, c|_{0.05} = 2.093)$$

$$2.53 \pm 2.47$$

