# Applied Machine Learning (AML)

Model Selection

Oisin Mac Aodha ● Siddharth N.

---

**Direct Comparison**

---

## Comparing Evaluation Measures

| email | true | pred (A) | pred (B) |
|---|---|---|---|
| "send us your password" | + | + | + |
| "send us review" | − | + | − |
| "review your account" | − | − | + |
| "review us" | + | − | − |
| "send your password" | + | + | + |
| "send us your account" | + | + | − |
| ⋮ | | | |

| | Naive Bayes (A) | Logistic Regression (B) |
|---|---|---|
| Acc | 72.6% | 84.5% |
| $\kappa$ | 54.1% | 66.2% |
| F1-score | 85.6% | 89.1% |
| ROC AUC | 48.4% | 55.7% |
| ⋮ | ⋮ | ⋮ |

Clearly, logistic regression (B) has higher scores than naive Bayes (A)!
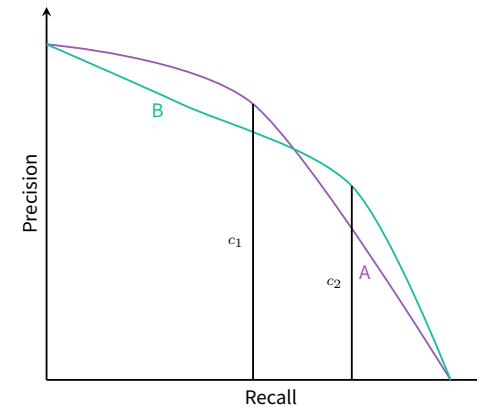
Should we choose B over A?    maybe?

---

## Comparing Point Estimates

$$\mathcal{D} = \{\mathcal{D}_{\text{train}}, \mathcal{D}_{\text{test}}\} \qquad \mathcal{D}_{\text{train}} \cap \mathcal{D}_{\text{test}} = \emptyset$$

| | Naive Bayes (A) | | Logistic Regression (B) |
|---|---|---|---|
| Acc | 72.6% | < | 84.5% |
| $\kappa$ | 54.1% | < | 66.2% |
| F1-score | 85.6% | < | 89.1% |
| ROC AUC | 48.4% | < | 55.7% |
| ⋮ | ⋮ | | ⋮ |

# Comparing Point Estimates

$$\mathcal{D} = \{\mathcal{D}'_{\text{train}}, \mathcal{D}'_{\text{test}}\} \qquad \mathcal{D}'_{\text{train}} \cap \mathcal{D}'_{\text{test}} = \emptyset$$

|  | Naive Bayes (A) | | Logistic Regression (B) |
|---|---|---|---|
| Acc | 79.3% | > | 78.1% |
| $\kappa$ | 61.9% | > | 60.3% |
| F1-score | 86.1% | > | 82.4% |
| ROC AUC | 50.1% | < | 50.4% |
| ⋮ | ⋮ | | ⋮ |

Point estimates can be susceptible to many kinds of random effects!

# Comparison with Tradeoff



## AUC of Precision-Recall

- Which model is better?
- Choice can depend on trade-off
  - lower recall, higher precision ($c_1$): $A > B$
  - lower precision, higher recall ($c_2$): $B > A$
- Random effects (e.g. data split) can make comparison hard

# Embracing Uncertainty

## Variation in error

- Dataset partitioning (e.g. cross validation)

$$\{\mathcal{D}^1_{\text{train}}, \mathcal{D}^1_{\text{test}}\}, \{\mathcal{D}^2_{\text{train}}, \mathcal{D}^2_{\text{test}}\}, \dots, \{\mathcal{D}^K_{\text{train}}, \mathcal{D}^K_{\text{test}}\}$$
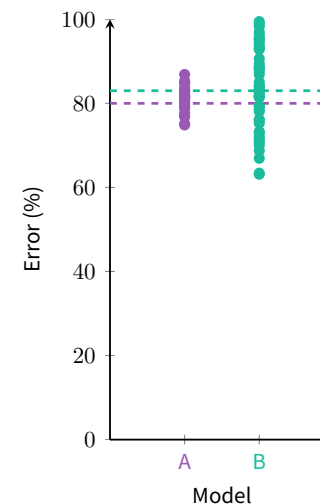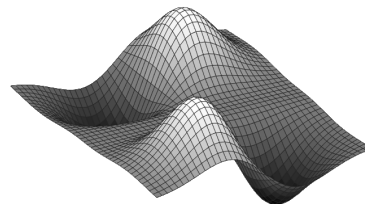
  A> B          A> B          …          B> A

- Model (e.g. stochastic linear regression)

$$y_i = w_0 + w_1 x_i + \epsilon_i \quad \epsilon_i \sim \mathcal{N}(0, 1)$$

- Learning algorithm (e.g. SGD)
  - initialisation effects
  - local minima

## Comparing Distributions

- Compute the difference in *mean* error
  - what difference is enough to decide B> A?
  - does the spread / variance affect this choice?

- Difficult to provide a general approach to say one model is "better" than another

- Weaker, but feasible, approach:

  How likely is it that the observed disparities are due to chance?

## Statistical Tests

---

## 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 $x_1, \ldots, x_N$

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

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

*Bessel's correction
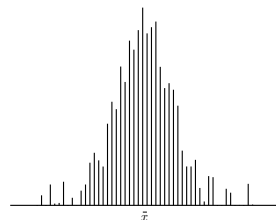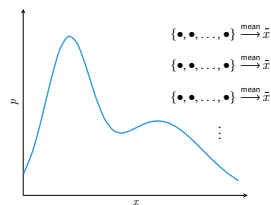
---

## Preliminaries

### Central Limit Theorem (CLT)

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

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

### Assume

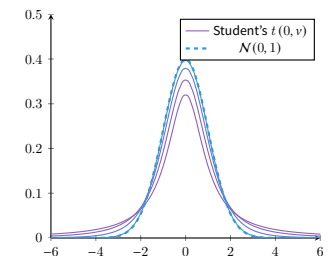- population $\mu$ known
- population $\sigma^2$ known

---

## Preliminaries

### Student's-$t$ distribution

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

### Assume

- population $\mu$ known
- population $\sigma^2$ *unknown*
- estimate sample variance $s^2 = \frac{1}{N-1} \sum_{i=1}^{N} (x_i - \bar{x}_N)^2$

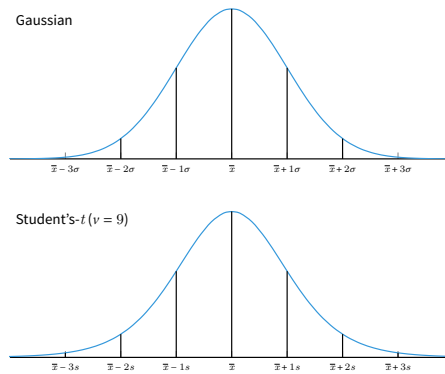$$t = \frac{\bar{x} - \mu}{s/\sqrt{N}}, \quad v = N - 1$$



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

# Statistical Testing: A Sketch

- Examine the *mean* of a set of samples
  e.g. difference in classification errors
- Why? — tendency towards Gaussian
- For some assumptions about the
  *population*: mean, variance (?)

  How likely is this observed sample mean
  value to have arisen by chance?

A common framework to evaluate
chance occurrence.

Gaussian

Student's-$t$ $(\nu = 9)$

**Statistical Tests**

**Hypothesis Testing**

---

# Hypothesis Testing

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

Decision (Retain)

|  |  | $H_0$ | $H_1$ |
|---|---|---|---|
| True | $H_0$ | ✓ | Type II |
|  | $H_1$ | Type I | ✓ |

### Null Hypothesis: $H_0$

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

### Alternative Hypothesis: $H_1$

- States the assumption believed to be True
- Evaluate if evidence supports assumption
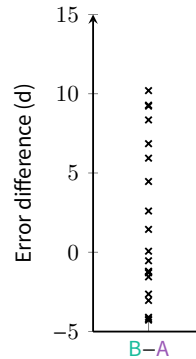- Always evaluates (strict) *in*equality $(\neq, >, <)$

---

# Hypothesis Testing: Variants

- Test type
  - $z$-test: Gaussian distribution
  - $t$-test: Student's $t$ distribution

- One or Two sided
  - **One:** $H_0 : \mu^A - \mu^B \leq 0$    $H_1 : \mu^A - \mu^B > 0$    (directional)
  - **Two:** $H_0 : \mu^A - \mu^B = 0$    $H_1 : \mu^A - \mu^B \neq 0$    (not directional)

- Test Statistic
  - **One-Sample:** compare sample to population with known characteristics
  - **Two-Sample:** compare two samples; typically experiment vs. control (e.g. vaccines)
  - **Paired:** one-sample test on *difference* between samples

# Example: Hypothesis Testing for Models

## Generating Variation

| Data Split | A | B | $d$ |
|---|---|---|---|
| $\{\mathcal{D}_{\text{train}}^1, \mathcal{D}_{\text{test}}^1\}$ | $\ell_1^A$ | $\ell_1^B$ | $\ell_1^B - \ell_1^A$ |
| $\{\mathcal{D}_{\text{train}}^2, \mathcal{D}_{\text{test}}^2\}$ | $\ell_2^A$ | $\ell_2^B$ | $\ell_2^B - \ell_2^A$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $\{\mathcal{D}_{\text{train}}^N, \mathcal{D}_{\text{test}}^N\}$ | $\ell_N^A$ | $\ell_N^B$ | $\ell_N^B - \ell_N^A$ |

---

# Example: Hypothesis Testing for Models



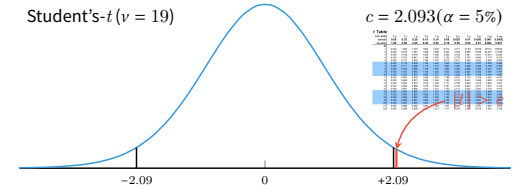| | |
|---|---|
| | −0.61 |
| | 6.44 |
| | −4.83 |
| | 7.43 |
| | −3.04 |
| | −4.75 |
| | 6.47 |
| | 9.62 |
| | 9.70 |
| | 2.51 |
| | −3.18 |
| | 8.23 |
| | 10.24 |
| | −4.59 |
| | 6.41 |
| | 2.08 |
| | −2.49 |
| | 2.72 |
| | −0.24 |
| | 2.56 |
| B−A | $d$ |

## Hypotheses

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

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

$$\bar{d} = \frac{1}{N}\sum_{i=1}^N d_i \qquad s^2 = \frac{1}{N-1}\sum_{i=1}^N (d_i - \bar{d})^2 \qquad t = \frac{\bar{d}-0}{s/\sqrt{N}}$$

$$= 2.53 \qquad\qquad = 27.78 \qquad\qquad = 2.14$$

Student's-$t\,(\nu = 19)$ $\qquad\qquad c = 2.093(\alpha = 5\%)$

---

# Hypothesis Testing: Caveats

- Rejecting $H_0$ does not imply 100% sure $H_0$ is False
- 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)

$$\bar{d} \pm c|_{\alpha/2} \cdot \frac{s}{\sqrt{N}}$$

which, for our example would be

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

$$2.53 \pm 2.47$$

---

# Cross Validation for Variation: Caveat

- Recall that CLT requires the samples to be independent
- Simple cross-validation can violate that independence (overlap in $\mathcal{D}_{\text{train}}$!)

| Data Split | A | B | $d$ |
|---|---|---|---|
| $\{\mathcal{D}_{\text{train}}^1, \mathcal{D}_{\text{test}}^1\}$ | $\ell_1^A$ | $\ell_1^B$ | $\ell_1^B - \ell_1^A$ |
| $\{\mathcal{D}_{\text{train}}^2, \mathcal{D}_{\text{test}}^2\}$ | $\ell_2^A$ | $\ell_2^B$ | $\ell_2^B - \ell_2^A$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |

- **Solutions:**
  - 5x2 Cross Validation [1]
  - Adjust standard deviation to account for imbalance [2]
  - …and many more (ANOVA, Non-parametric tests, etc.)!

1. T. G. Dietterich, Approximate Statistical Tests for Comparing Supervised Classification Learning Algorithms, 1998
2. C. Nadeau & Y. Bengio, Inference for the Generalization Error, 2003

# Summary

### Key

Being able to compare models and experiments is both a science and an art!

Most important aspect is to think what sources of variability affects results, and how large their effects are likely to be.

- Some measures incorporate context; use it! (P-R, ROC)
- For when statistical tests are required (not always!)
  - ensure your assumptions on the model / data are clearly stated
  - ensure assumptions of the test are met
- Performance on error measures not all—speed, use of resources, and ease of implementation can, and should, affect preference!

THE UNIVERSITY of EDINBURGH
**informatics**