

RIT - GCCIS SYLLABUS  
ISTE.780.01 DATA-DRIVEN KNOWLEDGE DISCOVERY  
SUMMER 2021 (TERM 2208) DRAFT OF MAY 30, 2021

DETAILS

Important note: The information presented in this syllabus is subject to expansion, contraction, change, or stasis during the semester. In case of conflict between versions, the copy on myCourses takes precedence.

**Course Number.** 80361

**Prerequisites.**

- ISTE-600 (Analytical Thinking) or equivalent
- PSYC-640 (Statistics) or equivalent

**Time.** 1000–1150

**Place.** GOL-2650 / Online

**Dates.** 26 May–4 Aug 2021

**Final Exam.** Online. Time To Be Announced.

**Instructor.** Mick McQuaid

**Email.** [mjmics@rit.edu](mailto:mjmics@rit.edu)

**Office.** 70-2675

**Office Hours.** by appointment at <https://rit.zoom.us/my/mjmics>

DESCRIPTION

Rapidly expanding collections of data from all areas of society are becoming available in digital form. Computer-based methods are available to facilitate discovering new information and knowledge that is embedded in these collections of data. This course provides students with an introduction to

the use of these data analytic methods, especially statistical learning approaches, within the context of the data-driven knowledge discovery process.

## MATERIALS

Each student will need a computing environment capable of running the R programming language. A minimal Linux, Mac, or Windows computer should suffice. R will be taught as part of this course. Students should have proficiency with some full-featured text editor. A network connection capable of streaming video is essential for this online course.

The textbook for this course is *An Introduction to Statistical Learning with Applications in R*, James et al. (2013). This book is available free online and in Wallace Library.

## SCHEDULE

Arabic numerals refer to days. The course runs for twenty days over eleven weeks (Roman numerals).

### Week I

1. May 27. The textbook — Chapter 1 (pp. 1–14) — Supervised and unsupervised learning — Prediction, inference, and description — Linear and non-linear models — Linear algebra — Vectors and matrices —  $n$  and  $p$  — The R language — Object orientation — The virtual machine — Examples — Grading — Individual work

### Week II

2. Jun 1. Chapter 2 (pp. 15–28) — Input and output, predictor and response, independent and dependent variable, features and target,  $X_i$  and  $Y$  — A basic model  $Y = f(X) + \epsilon$  — Prediction, inference, and description

— Irreducible and reducible error — Notation for expected value of error  $E(Y - \hat{Y})^2$  and variance of error — Examples — Estimating  $f$  as  $\hat{f}$  — Training data — Parametric methods — Estimating parameters — Non-parametric methods — Accuracy, flexibility, and interpretability — Supervised and unsupervised learning — Regression and classification — Exercises 2.4.1, 2.4.2

3. Jun 3. Chapter 2 (pp. 29–42) — Measuring the quality of fit — Mean squared error — Training and test data — Flexibility or degrees of freedom — Overfitting — Bias, error and variance — Classification — Indicator variables — Training error rate — Test error rate — Bayes classifier — Conditional probability — Bayes decision boundary — Bayes error rate —  $K$ -nearest neighbors or KNN — Exercise 2.4.9 — Team formation deadline

### Week III

4. Jun 8. Chapter 2 (pp. 42–57) — Introduction to R — Basic R commands — R graphics — Indexing data — Loading data — Graphical summaries of data — Numerical summaries of data — Exercise 2.4.8
5. Jun 10. Chapter 3 (pp. 59–101) — Linear Regression — Estimating coefficients — Assessing the accuracy of the coefficients — Assessing the accuracy of the model — Multiple linear regression — Estimating multiple regression coefficients — Questions about multiple linear regression — Qualitative predictors — Extensions of the linear model — Potential problems with the linear model — Exercise 3.7.8 — Milestone 1: Proposal

### Week IV

6. Jun 15. Chapter 3 (pp. 102–126) — An extended example of multiple linear regression — Lab on linear regression — Libraries in R for linear regression — Simple linear regression in R — Multiple linear regression in R — Interaction terms in R — Nonlinear transformations of predictors — Qualitative predictors — Writing functions — Exercise 3.7.9
7. Jun 17. Chapter 4 (pp. 128–154) — Classification — Overview — Alternative to linear regression — Logistic regression — Logistic model — Estimating logistic regression coefficients — Making predictions — Multiple logistic regression — Logistic regression with more than two response classes — Linear discriminant analysis — Using Bayes' theorem for classification — Linear discriminant analysis in one dimension — Linear discriminant analysis in more than one dimension — Quadratic discriminant analysis — Comparing classification methods — Exercise 4.7.I

## Week V

8. Jun 22. Chapter 4 (pp. 154–174) — Lab on logistic regression, LDA, QDA, and KNN in R — Stock market data — Logistic regression in R — Linear discriminant analysis in R — Quadratic discriminant analysis in R —  $K$ -nearest neighbors in R — Application to caravan insurance data — Exercise 4.7.II
9. Jun 24. Chapter 5 (pp. 175–190) — Resampling — Cross-validation — Validation set approach — Leave-one-out cross-validation —  $k$ -fold cross-validation — Bias variance trade-off — Cross-validation on classification problems — Bootstrap — Exercise 5.4.5

## Week VI

10. Jun 29. Chapter 5 (pp. 190–202) — Lab: Cross-validation and bootstrap in R — Validation set approach in R — Leave-one-out cross-validation in R —  $k$ -fold cross-validation in R — Bootstrap in R — Exercise 5.4.6
11. Jul 1. Chapter 6 (pp. 203–244) — Linear model selection — Subset selection — Stepwise selection — Choosing the optimal model — Shrinkage methods — Ridge regression — Lasso — Selecting the tuning parameter — Dimension reduction methods — Principal components regression — Partial least squares — Considerations in high dimensions — High-dimensional data — What can go wrong in high dimensions — Regression in high dimensions — Interpreting results in high dimensions — Exercise 6.8.8 — Milestone 2: Data summary / visualization

## Week VII

12. Jul 6. Chapter 6 (pp. 244–264) — Lab: Subset selection in methods in R — Forward and backward stepwise selection in R — Choosing among models using the validation set approach and cross-validation in R — Lab: Ridge regression and the lasso — Ridge regression in R — Lasso in R — Lab: Principal components regression and partial least squares in R — Principal components regression in R — Partial least squares in R — Exercise 6.8.10
13. Jul 8. Chapter 7 (pp. 265–280) — Beyond linearity — Polynomial regression — Step functions — Basis functions — Regression splines — Piecewise polynomials — Constraints and splines — Spline basis representation —

Choosing the number and locations of knots — Comparison between regression splines and polynomial regression — Smoothing splines — Overview of smoothing splines — Choosing the smoothing parameter  $\lambda$  — Exercise 7.9.6

## Week VIII

14. Jul 13. Chapter 7 (pp. 280–301) — Local regression — Generalized additive models — GAMs for regression problems — GAMs for classification problems — Lab: non-linear modeling — Polynomial regression and step functions in R — Splines in R — GAMs in R — Exercise 7.9.7
15. Jul 15. Chapter 8 (pp. 303–323) — Tree-based methods — Basics of decision trees — Regression trees — Classification trees — Trees versus linear models — Advantages and disadvantages of trees — Bagging, random forests, boosting — Bagging — Random forests — Boosting — Exercise 8.4.3 — Milestone 3: Algorithm testing

## Week IX

16. Jul 20. Chapter 8 (pp. 323–336) — Lab: Decision trees — Fitting classification trees in R — Fitting regression trees in R — Bagging and random forests in R — Boosting in R — Exercise 8.4.8
17. Jul 22. Chapter 9 (pp. 337–354) — Support vector machines — Maximal margin classifier overview — Definition of a hyperplane — Classification using a separating hyperplane — Maximal margin classifier — Constructing the maximal margin classifier — Non-separable

case — Support vector classifiers — Overview of the support vector classifier — Details of the support vector classifier — Support vector machine overview — Classification with non-linear decision boundaries — Support vector machines — Application to the heart disease data — Exercise 9.7.4

## Week X

18. Jul 27. Chapter 9 (pp. 355–372) — Support vector machines with more than two classes — One-versus-one classification — One-versus-all classification — Lab: Support vector machines — Support vector classifier in R — Support vector machine in R — ROC curves in R — SVM with multiple classes in R — Application to gene expression data — Exercise 9.7.5
19. Jul 29. Chapter 10 (pp. 373–400) — Unsupervised learning — Challenge of unsupervised learning — Principal components analysis — Defining principal components — Alternative interpretation of principal components — Details of principal component analysis — Other uses for principal components — Clustering methods —  $K$ -means clustering — Hierarchical clustering — Practical issues in clustering — Exercise 10.7.9 — Milestone 4: Core algorithm tuning

## Week XI

20. Aug 3. Chapter 10 (pp. 401–413) — Lab: Principal components analysis — Lab: Clustering —  $K$ -means clustering in R — Hierarchical clustering in R — Lab: NCI60 data example — Principal components analysis on the NCI60 data — Clustering the observations of the NCI60 data — Exercise 10.7.10

Aug 4. Milestone 5: Final project paper

## GRADING

The grading scale used along with the grade components follow.

- A  $\geq 90.0\%$
- B  $\geq 80.0\%$  &  $< 90.0\%$
- C  $\geq 70.0\%$  &  $< 80.0\%$
- D  $\geq 60.0\%$  &  $< 70.0\%$
- F  $< 60.0\%$

The course grade is composed of fifteen percent for each of five milestones and twenty-five percent for the final exam for a total of 100 percent.

## POLICIES

The following are brief statements of policy that are, in many places, expanded at the URLs provided. You are bound by these policies and any protest that you did not read the extended versions at the provided links will not be heeded. Your familiarity with the following policies, dates, and parameters will be assumed in this course.

**Last day of 7-day add/drop period.** Wednesday 2 Jun 2021

**Last day to withdraw with W.** 22 July 2021

**myCourses.** All project assignments, lecture notes, and other distributable course materials will be available via myCourses. Except where otherwise indicated, all student project assignments will be submitted via myCourses dropboxes.

**Grade Challenges.** School of Information policy states that a student has one semester to challenge any grade. After that, grades cannot be challenged.



**Late Work.** Any work not submitted by the final due date receives a grade of zero unless arrangements are made previous to the initial due date.

**Extra Credit.** No extra credit is available in this course.

**Accommodations.** If you have a “Notice of Accommodation”, you must provide your instructor with a copy of it within 1 week of starting this course. You must follow all the rules of the relevant office.

**Academic Dishonesty.** The policy on dishonesty is simple: Anyone caught cheating receives an “F” as a course grade, is removed from the section and a letter detailing the incident is placed into his or her folder. Any student accused of cheating should realize that the evidence has already been verified by other faculty members and will withstand an appeal. Additionally, please review the institute policy at [http://www.rit.edu/studentaffairs/studentconduct/rr\\_academicdishonesty.php](http://www.rit.edu/studentaffairs/studentconduct/rr_academicdishonesty.php)

**Acceptable Use.** We are bound by the following Acceptable Computer Use policy at <http://www.rit.edu/academicaffairs/policiesmanual/sectionC/C82.html>

**Student Responsibilities.** Please review the general student responsibilities as outlined at <http://www.rit.edu/~301www/rr.php3>

**Policy on Reporting Incidents of Discrimination and Harassment.** RIT is committed to providing a safe learning environment, free of harassment and discrimination as articulated in our university policies located on our governance website. RIT’s policies *require faculty to share information* about incidents of gender based discrimination and harassment with RIT’s Title IX coordinator or deputy coordinators, regardless whether the incidents are stated to them in person or shared by students as part of their coursework. RIT

Governance website: <https://www.rit.edu/academicaffairs/policiesmanual/policies/governance>

If you have a concern related to gender-based discrimination and/or harassment and prefer to have a *confidential* discussion, assistance is available from one of RIT's confidential resources on campus:

1. The Center for Women & Gender: Campus Center Room 1760; 585-475-7464; CARES (available 24 hours/7 days a week) Call or text 585-295-3533.
2. RIT Student Health Center – August Health Center/1st floor; 585-475-2255.
3. RIT Counseling Center - August Health Center /2nd floor - 2100; 585-475-2261.
4. The Ombuds Office – Student Auxiliary Union/Room III4; 585-475-7200 or 585-475-2876.
5. The Center for Religious Life – Schmitt Interfaith Center / Rm 1400; 585-475-2137.
6. NTID Counseling & Academic Advising Services – 2nd Floor Lyndon B. Johnson; 585-475-6468 (v), 585-286-4070 (vp).

**RIT Resilience.** Success in this course depends heavily on your personal health and wellbeing. Recognize that stress is an expected part of the college experience, and it often can be compounded by unexpected setbacks or life changes outside the classroom. Your other instructors and I strongly encourage you to reframe challenges as an unavoidable pathway to success. Reflect on your role in taking care of yourself throughout the term, before the demands of exams and projects reach their peak. Please feel free to reach out to me about any difficulty you may be having that may impact your performance in this course as soon as it occurs and before it becomes unmanageable. In addition to your academic advisor,

I strongly encourage you to contact the many other support services on campus that stand ready to assist you.