

Instructions for the SI630 Project Proposal

Version 1.0

Put your name here

1 Introduction

The course project is intended to provide an opportunity for students to dive deeper into one problem or topic of their choice and write a very small scale study on the topic. Projects typically take two forms: (1) the student has some data, problem, or algorithm in mind and proposes a study to investigate these or (2) students pick an existing NLP task and try a new approach to solving it. Tasks for the latter are detailed more in Section ???. In both cases, projects should be *feasible* for completing in the available time frame. The latter part of the course has a lighter workload to allow more time for working on the project, but we want to ensure that students pick projects that help them learn real, practical skills in NLP without being trivial. Ideally, your course project is a chance to develop something you can show off to future employers or could serve as a pilot study for a full research project.

Learning goals for the course project are as follows.

1. Learn advanced NLP develop skills through practice
2. Gain specialized knowledge in a particular topic or problem
3. Provide an end-to-end experience from data to results in answering a question using NLP
4. Practice forming a research question and designing a series of experiments to answer that question
5. Learn about new NLP methods through literature search
6. Learn about a particular topic or problem through literature search

7. Practice preparing high-quality technical reports in typesetting (latex)
8. Produce a concrete artifact that can be shown to interested parties (employers).

Please remember that the 630 instructors are here for you and will gladly offer suggestions and advice on projects. We want your projects to succeed, to be fun to work on, and to spark your intellectual curiosity!

2 Formatting

For this project, you're expected to use this L^AT_EX template. You're welcome to copy this template directly off of Overleaf as well using the ACL 2021 template at this URL: <https://www.overleaf.com/latex/templates/instructions-for-acl-ijcnlp-2021-proceedings-mhxffkjdwymb>, which hopefully has enough examples of how things can be written to get you started. If you're having any issues getting L^AT_EX to do what you want please feel free to ask the instructors or know that there's a great resource on WikiBooks <https://en.wikibooks.org/wiki/LaTeX> and a whole StackExchange site dedicated to answering questions <https://tex.stackexchange.com/>. L^AT_EX is a common method for writing technical documents so we are using it for 630 to help get you started on using it for your career. It also is pretty awesome for citation management.

3 What to Cover

Your proposal is expected to address the following points in a coherent document that reads like a proposal—not a list of bullet points. We typically expect that proposals are written in a coherent structure with sections denoting coherent elements. Often, many students write their proposal

in a way that it can be re-used as the starting point for the Project Update and later in the Final Report. The goal of the proposal is to help you get started on that writing process.

Project Goals (0.5 points) For the proposal, you should have a rough draft of the introduction that clearly states what are the goals of the project is and provides some broader context for why these are important goals. You should also include a statement on why solving this problem matters—who would care if you solved it and what effect would solving it have? The SI630 projects will all be made into blog posts which typically attract a broad audience. As an eye to the future, in the proposal, think about who you would want to see your work and why your results would matter to different groups of people.

NLP Task Definition (1 points) The section describes what specific problem or NLP Task you plan to solve and goes into much more detail than what’s specified in the Project Goals in the Introduction. You can also add details about what your problem is not to help guide the reader’s expectations. For the proposal, describe very clearly what problem you will solve. It helps to be specific about what kind of input your system will use and what specifically it will produce as output.

Data (0.5 points) This section goes into detail what data you will use to train and evaluate a model for your particular NLP Task. Ideally, you should already have the data on hand or spend a few minutes getting it. Projects are *much* more successful when most of the time is spent solving the problem rather than searching for the right data. If you don’t have data at proposal time, please be very specific on how you plan to get it. We might be able to recommend something but we encourage you to come to talk to us during office hours or after class. If you have the data, you should include a few examples and can also include some very rough statistics (e.g., how many instances you have). Tables and figures are very useful for showing examples. Please make sure the text is legible!¹

Since not everyone has a burning question they’re dying to answer with NLP, we’ve included a few potential NLP problems people could work

¹Copying images of tables/data in as figures often results in illegible text. If you’re having trouble getting figures to work, come talk to us!

on in Table 1. For these tasks, the problem, data, and evaluation criteria are already provided—though you should be sure to describe them here as a part of your proposal. **Important note:** projects *cannot* be derived from any Kaggle competition or similar type of shared task, without prior approval. Essentially, if the instructors search online for your dataset and problem and can find a fully-worked example on a blog/github, the idea can’t be used. Submissions using this kind of data will be turned back with a grade of zero and you’ll have to resubmit.

If you’re stuck for looking for data try using one of the many dataset search engines from [Microsoft](#) or [Google](#). In general, we recommend finding a topic area that you’re passionate about, thinking of a few questions that involve text from that area, and then looking for datasets. If you’re stuck, ask on Piazza and we can collectively try to find something.

Related Work (0.5 points) You should have at least *three* papers related to your current problem and a few sentences describing what they did to solve the problem. The related work section should describe how other people have thought about the problem you’re working on. How did they approach it? What makes their problem different from yours? Why do you think your approach will be better?

Be sure to cite the papers using $\text{BIB}\text{T}_{\text{E}}\text{X}$ (see the references.bib file this overleaf project). If you’re not sure which papers to cite, try searching for your topic and keywords on Google Scholar or Semantic Scholar; their “related papers” functionality can help you find interesting papers for ideas.

Proposed Method (0.5 points) The proposal should include a high-level description of how you plan to solve the the problem. You don’t need to have specific details of implementations but we encourage you to start thinking about it now. You might take some insight from relate works to see general strategies. As a part of this section, please describe *why* you chose that approach and why you think it will successful.

Evaluation Plan (0.5 point) You *must* include an specific evaluation metric as a part of the proposal. In essence, you need to think about how you will measure how good your NLP system is at your particular task or at solving your problem. Manual analysis of output can be helpful, but you

Task	Reference Paper	Website
Natural Language Inference	(Bowman et al., 2015)	https://nlp.stanford.edu/projects/snli/
Identifying which parts of a message are toxic		https://sites.google.com/view/toxicspans
Extracting scientific findings from papers		https://ncg-task.github.io
Reasoning about facts in tabular data		https://sites.google.com/view/sem-tab-facts

Table 1: Examples of NLP tasks that you could choose to work on for your project. The last three tasks are very recent tasks that have details and data on the website but no papers yet. You can work on something cutting edge!

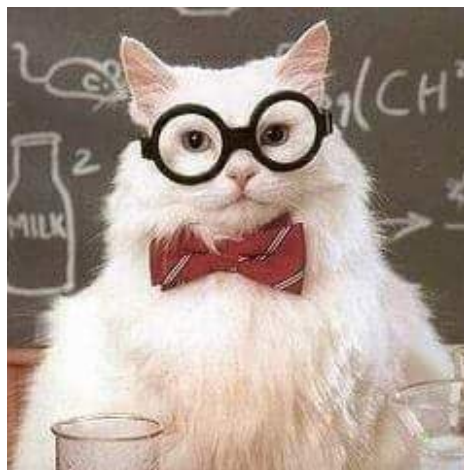


Figure 1: Eventually, you'll have cool figures to show here.

need to know quantitatively how to measure goodness.² If you're not sure about a good evaluation, please post of Piazza and we can work it out there.

The second part of the evaluation plan is what baselines you'll compare your system against. One baseline should be random performance. A second baseline should be something reasonable that doesn't require much knowledge or learning. For example, if you're doing a classification, always choosing the most frequent class is a useful baseline. A baseline is essential here because it helps you figure out what is the *simplest* way to solve your task. Baselines provide a useful comparison for putting your model in context—several students in past semesters have been surprised by how well baselines performed.

Work Plan (0.5 points) You'll end with a plan for how you'll accomplish your project. We hope this section can help you think at a high-level about which tasks are necessary to get to the point where you have a working model. Of course, plans

are often made and then changed when new information or challenges emerge, so we won't hold you to this. However, the act of writing a plan can greatly help you figure out how think about the process and, in general, projects that are proposed with more concrete work plans tend to be more successful.

Multi-person Team Justification If you have more than one person on your project, you should justify why the work requires the number of people you have. In addition, you should provide a concrete explanation of what each team member is expected to do, keeping in mind that all team members need to be doing *some* kind of NLP for the project. An n person project should have $n * 1.25$ peoples worth of work.

References

Samuel R Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. 2015. A large annotated corpus for learning natural language inference. *arXiv preprint arXiv:1508.05326*.

²As a side note, projects involving topic modeling on a corpus usually have a harder time coming up with a precise metric (i.e., how good is one topic model versus another?), so I would avoid those unless you have an extrinsic tasks that you can test on (i.e., topic modeling is used as a *method* to facilitate solving some other kind of task).