

Pr. 1.

Let π be an N -dimensional vector having nonnegative elements that sum to one.

- Give a **transition matrix** $P \neq I$ of a **Markov chain** having π as an **equilibrium distribution**.
- For that transition matrix P , discuss whether the given π is the *unique* equilibrium distribution.
- For that transition matrix P , discuss whether the Markov chain **power iteration** $\pi_{k+1} = P\pi_k$ is guaranteed to converge to π for any initial probability distribution π_1 .
- Optional. Discuss whether there are additional answers to (a).

Pr. 2.

Consider the **Markov chain** that has **transition matrix** $P = G_N + (1-p)(e_3 - e_2)e_1'$ for $N \geq 3$, where G_N denotes the circulant generator matrix, and $p \in (0, 1)$. Determine the equilibrium distribution π_* of this Markov chain and discuss whether it is unique. This is a generalization of the Markov chain shown in the demo.

Hint. If needed, solve it numerically for a few small values of $N \geq 3$ to form a conjecture that you can then verify analytically. Your solution must be general for $N \geq 3$ and $p \in (0, 1)$. For reference:

$$P = \begin{bmatrix} 0 & 0 & \dots & 0 & 1 \\ p & 0 & 0 & \dots & 0 \\ 1-p & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \dots & 1 & 0 \end{bmatrix}.$$

Pr. 3.

The mythical Michigan Wolverine eats something once a day, either cheese, grapes, or lettuce. Its daily dietary habits follow the following rules:

- If it ate cheese today, tomorrow it will eat lettuce or grapes with equal probability.
- If it ate grapes today, tomorrow it will eat grapes with probability 1/10, cheese with probability 4/10, and lettuce with probability 5/10.
- If it ate lettuce today, it will not eat lettuce again tomorrow, but it will eat grapes with probability 4/10 or cheese with probability 6/10.



- What is the **long term average probability** of the Wolverine eating cheese, grapes, or lettuce on some day?
- Is your answer the unique solution? Explain why or why not.

Pr. 4.

An alternative to the logistic loss function for **binary classifier design** is the ℓ_2 -hinge **loss function**:

$$\psi(t) \triangleq (\max(1-t, 0))^2.$$

- Show that this loss function has a Lipschitz continuous derivative and determine its Lipschitz constant \mathcal{L}_ψ .
- Now consider **logistic regression** using this loss function for binary classifier design using the optimization problem

$$\arg \min_{\mathbf{x}} f(\mathbf{x}), \quad f(\mathbf{x}) = \sum_{m=1}^M \psi(y_m \langle \mathbf{x}, \mathbf{v}_m \rangle) + \beta \frac{1}{2} \|\mathbf{x}\|_2^2.$$

Determine a Lipschitz constant $\mathcal{L}_{\nabla f}$ for the gradient of the overall cost function $f(\mathbf{x})$, where $\mathbf{x}, \mathbf{v}_m \in \mathbb{R}^N$ and $y_m = \pm 1$.

Pr. 5.

- (a) Derive the optimal **step size** α_* for classical **gradient descent** of the **least-squares** cost function $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2$, in terms of the singular values of the (full column rank) matrix \mathbf{A} . Here, optimal means $\rho(\mathbf{I} - \alpha\mathbf{A}'\mathbf{A})$ is as small as possible.
- (b) For that optimal step size α_* , determine the **spectral radius** $\rho(\mathbf{I} - \alpha_*\mathbf{A}'\mathbf{A})$ that governs the **convergence rate** of the sequence $\{\mathbf{x}_k\}$ produced by the GD method.
Express your answer in terms of the **condition number** of $\mathbf{A}'\mathbf{A}$.
-

Pr. 6.

Practical preconditioners are often built on approximations of a cost function's Hessian inverse. The initial design of such preconditioners often does not ensure convergence of PGD, so a step size α is usually needed. Specifically, for a LS problem we might have $\mathbf{P}_0 \approx (\mathbf{A}'\mathbf{A})^{-1}$ and then let $\mathbf{P} = \alpha\mathbf{P}_0$, where \mathbf{P}_0 is positive definite.

- (a) Determine the optimal (fixed) **step size** α_* that makes the **preconditioned gradient descent** (PGD) iteration $\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha\mathbf{P}_0\mathbf{A}'(\mathbf{A}\mathbf{x}_k - \mathbf{y})$ have the fastest possible **convergence rate**.
Hint. The answer depends on both \mathbf{A} and the preconditioner \mathbf{P}_0 .
- (b) Simplify your expression for α_* as much as possible in the special case where we use the ideal preconditioner $\mathbf{P}_0 = (\mathbf{A}'\mathbf{A})^{-1}$. Does your simple expression (a number) make sense?
-

Pr. 7.

On HW2 we reminded you that your graduate school experience should be more than just taking courses, and we encouraged you to write down some of your own non-course-related goals for this semester.

For 10 points, reflect on how it went and write down some of your thoughts about what you accomplished and did not complete of your original non-course goals. The graders will assign credit for anything that looks sincere. Prof. Fessler will read these after finals are graded.

Submit your reflection to this google form <https://forms.gle/L1GsfCLXW6ExwKYs8> and also upload it to gradescope. Thanks! We encourage you to keep this practice of planning and self reflection in future semesters.

Pr. 8.

Please complete the online course evaluation for *both* the lecture and your discussion section before the deadline given by UM. The system will give you a confirmation page that shows your name. Submit a screenshot of that confirmation (with your name) to gradescope to earn credit for this problem. Your evaluation is very important to us for improving the course in future terms. Your evaluation is entirely anonymous.

Pr. 9.**Logistic Regression Classification (Discussion Task)**

Download the `task-6-classify-logistic.ipynb` jupyter notebook file from Canvas under the discussion folder and follow all instructions to complete the task. You may work individually, but we recommend that you work in pairs or groups of three.

When you are finished, upload your solutions to gradescope. Note that the submission for the task is separate from the rest of the homework because the task allows you to submit as a group. Only upload one submission per group! Whoever uploads the group submission should add all group members in gradescope, using the “View or edit group” option on the right-hand sidebar after uploading a PDF and matching pages. Make sure to add all group members, because this is how they will receive credit.

As part of this task, you will need to complete the following Nesterov gradient descent function and submit it to the autograder *individually*.

In Julia, your file should be named `ngd.jl` and should contain the following function:

```
"""
    x, out = ngd(grad, x0 ; niter::Int = 200, L::Real = 0, fun = (x, iter) -> 0)

Nesterov's FGM (fast gradient method),
given the gradient of the cost function

# In:
- `grad` function that takes `x`
  and calculates the gradient of the cost function with respect to `x`
- `x0` an initial point

# Option:
- `niter` number of iterations
- `L` Lipschitz constant of the gradient of the cost function
- `fun` function of `(x, iter)` to evaluate every iteration

# Out:
- `x` guess of the minimizer after running `niter` iterations
- `out` `Vector` of `niter+1` evaluations of the `fun` function
"""
function ngd(grad::Function, x0::AbstractArray ;
    niter::Int = 200, L::Real = 0, fun::Function = (x, iter) -> 0)
```

Email your solution as an attachment to eeecs551@autograder.eecs.umich.edu.

Non-graded problem(s) below

(Solutions will be provided for self check; do not submit to gradescope.)

Pr. 10.**Low-rank matrix completion from noisy data**

In this problem you will use the **Schatten p-norm** regularizer for **matrix completion**, and show that it works better than the **nuclear norm** regularizer. In matrix completion, instead of being given a full matrix $\mathbf{Y} = \mathbf{X}_{\text{true}} + \boldsymbol{\varepsilon}$, where \mathbf{X}_{true} is low rank and $\boldsymbol{\varepsilon}$ is additive noise, we are given $\mathbf{Y} = \mathbf{M} \odot (\mathbf{X}_{\text{true}} + \boldsymbol{\varepsilon})$, where \mathbf{M} is a binary matrix (a sampling mask or pattern) that is 1 where we have data and 0 where we do not, and \odot denotes elementwise multiplication.

We want to estimate the matrix \mathbf{X}_{true} by solving the following optimization problem:

$$\hat{\mathbf{X}} = \arg \min_{\mathbf{X} \in \mathbb{C}^{M \times N}} \frac{1}{2} \|\mathbf{M} \odot (\mathbf{Y} - \mathbf{X})\|_{\text{F}}^2 + \beta R(\mathbf{X}),$$

where the regularizer $R(\mathbf{X})$ is the one associated with the Schatten p -norm for $p = 1/2$ as discussed in a previous HW problem.

- (a) Implement a **Julia** function that performs a specified number of iterations of FISTA to solve for $\hat{\mathbf{X}}$. As a guide for this, use code from the FISTA method in the notebook linked at this demo (as html) and here (as a notebook). That notebook uses the nuclear norm regularizer, whereas here you use the Schatten-based regularizer. Just like was done in the notebook, initialize with $\mathbf{X}_0 = \mathbf{Y}$, the “zero-filled” data.

You might want to debug this in a Jupyter notebook, but then extract your FISTA iteration as a function for submitting to the autograder.

In **Julia**, your file should be named `fista_schatten.jl` and should contain the following function:

```
"""
    Xh = fista_schatten(Y, M, reg::Real, niter::Int)

Perform matrix completion by using `niter` FISTA iterations
to seek the minimizer over `X` of
`1/2 || M .* (Y - X) ||^2 + reg R(x)`
where `R(X)` is the Schatten p-norm of `X` raised to the `p`th power,
for `p=1/2`, i.e., `R(X) = \sum_k (\sigma_k(X))^(1/2)`.

In:
- `Y` : `M × N` matrix (with zeros in missing data locations)
- `M` : `M × N` boolean matrix (`true` in sampled data locations)
- `reg` : regularization parameter
- `niter` : # of iterations

Output
- `Xh` : `M × N` estimate of `X` after `niter` FISTA iterations
"""
function fista_schatten(Y, M::AbstractMatrix{Bool}, reg::Real, niter::Int)
```

Email your solution as an attachment to `eeecs551@autograder.eecs.umich.edu`.

Hint. The solution requires only small changes to the FISTA code in the demo notebook.

- (b) Incorporate your `fista_schatten` code into the Jupyter notebook `09_lrnc_nuc` into the appropriate places at the end of the notebook, either by cut-and-paste or by `include("file.jl")` statements. (Be sure your FISTA method passes the autograder test first.) Apply your Schatten-based FISTA iteration to the “UM” data used in that notebook using $\beta = 120$ for the Schatten case.

Make images of the estimate $\hat{\mathbf{X}}$ and the error $\hat{\mathbf{X}} - \mathbf{X}_{\text{true}}$ for both the FISTA/SVST solution (with $\beta = 0.8$ per the notebook) and for your FISTA/Schatten solution (with $\beta = 120$ per the notebook). As configured in the notebook, show the NRMSE $\|\hat{\mathbf{X}} - \mathbf{X}_{\text{true}}\|_{\text{F}} / \|\mathbf{X}_{\text{true}}\|_{\text{F}} \cdot 100\%$ for the two methods.

- (c) Add a cell to the notebook to plot the cost function.
- (d) Discuss whether the cost function decreases monotonically, both empirically and theoretically.

- (e) Optional. Experiment with the regularization parameter β for one or both methods to see if you can make better images.

This is the final segment of the three-part matrix completion problem. You will have implemented a method using quite contemporary methods. Note that FISTA was derived for convex problems whereas for $p < 1$ the Schatten regularizer is a non-convex “quasi-norm” so convergence is (to my knowledge) an open question.

An interesting (but probably nontrivial) extension of this work would be to generalize **SURE** to select the regularization parameter β automatically.
