

Machine Learning and Large Scale Data Analysis

HW 1

Out: Tuesday, April 2, 2018

Due: Tuesday, April 9, 2018 (at 2:00 p.m.)

Please hand in this Homework in 5 files:

1. A pdf of the theoretical homework.
2. A pdf of your jupyter notebook for problem 4.
3. The ipynb file for problem 4.
4. A pdf of your jupyter notebook for problem 5.
5. The ipynb file for problem 5.

1. *Trouble with algorithms* (R points) read these articles

www.theguardian.com/science/2016/sep/01/how-algorithms-rule-our-working-lives

<https://www.nature.com/articles/d41586-018-05707-8>

and post a comment on

https://canvas.uchicago.edu/courses/21205/discussion_topics/167354

You can then 'like' comments of other people that you find interesting.

2. *Maximum likelihood* (10 points)

Maximum likelihood is a method to estimate parameters of a distribution from observations. Let $f(x, \theta), \theta \in \Theta$ be a family of distributions. Assume X_1, \dots, X_n are i.i.d samples from $f(x, \theta^*)$. For any value of $\theta \in \Theta$ the log-likelihood is

$$\ell(X_1, \dots, X_n; \theta) = \ell(\mathbf{X}, \theta) = \sum_{i=1}^n \log f(X_i; \theta).$$

This is minimized over Θ by solving the score equation

$$\frac{\partial \ell(\theta)}{\partial \theta} = 0.$$

The solution $\hat{\theta}$ is call the *maximum likelihood* estimate of θ . In most situations we will deal with this has a unique solution that is indeed the maximum.

- (a) Let X_1, \dots, X_n be independent draws from a Poisson distribution with mean λ . $P(X = k; \lambda) = e^{-\lambda} \lambda^k / k!$. Write the log-likelihood $\ell(\mathbf{X}, \lambda)$, derive the score equation and find $\hat{\lambda}$ that solves the score equation.

- (b) $X_1, \dots, X_n \in \mathbb{R}^d$ be i.i.d draws from the normal $N(\mu, \Sigma)$, with $\mu \in \mathbb{R}^d$ and Σ a positive definite $d \times d$ matrix. Write the score equation for μ and solve it (note: the solution does not depend on Σ).
- (c) Assume Σ is diagonal $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_d)$. Write the score equations for μ and $\sigma_1, \dots, \sigma_d$. Denote by $\hat{\mu}$ the maximum likelihood estimate that you obtained in part (b). Solve the score equations for $\hat{\sigma}_i, i = 1, \dots, d$.
- (d) Assume $\Sigma_0 \in \mathbb{R}^{d \times d}$ is positive definite and known and that $X_i \sim N(\mu, \alpha \Sigma_0)$, with $\theta = (\mu, \alpha)$ unknown. Write the score equation for α and solve for $\hat{\alpha}$.

3. Regression (10 points)

In linear regression, the fitted values are defined to be $\hat{\mathbf{y}} = X\hat{\beta}$ where $\hat{\mathbf{y}} = H\mathbf{y}$ and

$$H = X(X^T X)^{-1} X^T.$$

assuming $n > d$ and $X^T X$ is nonsingular. The matrix H is called the “hat matrix.” Define \mathcal{L} to be the set of vectors that can be obtained as linear combinations of the columns of X , which is an $n \times d$ matrix. Show that the hat matrix satisfies the following properties:

- (a) $\hat{\mathbf{y}} = H\mathbf{y} = X\hat{\beta}$ are the least squares estimates.
- (b) $HX = X$.
- (c) H is symmetric: $H = H^T$.
- (d) H is idempotent: $H^2 = H$.
- (e) $\hat{\mathbf{y}} = H\mathbf{y}$ is the projection of \mathbf{y} onto the column space \mathcal{L} .
- (f) $\text{rank}(X) = \text{tr}(H) = d$.

4. Singular value decomposition (10 points)

Let $X \in \mathbb{R}^{m \times n}$ have $\text{rank}(X) = r \leq \min(m, n)$ and let $X = U\Sigma V^T$ be the SVD of X where $U \in \mathbb{R}^{m \times m}$, $V \in \mathbb{R}^{n \times n}$ are orthogonal matrices and $\Sigma \in \mathbb{R}^{m \times n}$ is the diagonal matrix of singular values, with $\Sigma_{ii} = \sigma_i$, $i = 1, \dots, \min(m, n)$.

- (a) Show that the columns of U are eigenvectors of XX^T and the columns of V are eigenvectors of $X^T X$. Determine what are the corresponding eigenvalues in terms of $\sigma_1, \dots, \sigma_r$.
- (b) If u_1, \dots, u_m and v_1, \dots, v_n are the columns of U and V , show that $Xv_i = \sigma_i u_i$ and $X^T u_i = \sigma_i v_i$.
- (c) Express the Frobenius norm $\|X\|_F = \sqrt{\sum_{ij} X_{ij}^2}$ in terms of the singular values σ_i .

- (d) Express $|\det(X)|$ in terms of the singular values σ_i (Hint: you will need to use orthogonality to express $|\det(U)|$ and $|\det(V)|$).
- (e) Assuming $X^\top X$ is invertible, express the hat matrix $H = X(X^\top X)^{-1}X^\top$ of linear regression in terms of the SVD.
- (f) Let the $m \times n$ matrix $\Sigma^{(k)}$, $1 \leq k \leq r$ be the diagonal matrix $\Sigma_{ii}^{(k)} = \sigma_i$ for $i = 1, \dots, k$ (and zeros in the remaining entries). The matrix $U\Sigma^{(k)}V^\top$ is known as the rank- k approximation for X . Express the least square regression estimate obtained using the rank- k approximation instead of X .

5. Self-fulfilling prophecies (20 points)

We want to show what kind of issues can arise from clustering when it is applied to real people and affects their choices.

- (a) Simulate 1000 points from a bivariate normal distribution $N(\mu, \Sigma)$ with

$$\mu = 0, \Sigma = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}.$$

Plot the data. Pretend these two variables represent the features an app uses to characterizes people's preferences and subsequently cluster them into 3 clusters.

- (b) Using the python package `sklearn.cluster.KMeans` fit a 3 cluster model to the data. Plot the three cluster centers and color the data points according to the clusters they are assigned. This is the assignment the app has chosen for this population.
- (c) Now modify each point in the data to move 1% closer to its assigned cluster center. $x_i = .99 * x_i + .01c_i$, where c_i is the cluster center assigned to the i 'th data point. This corresponds to a tiny indirect effect of the choice of cluster on the features of the people in the sample. Now repeat the clustering on the modified data.
- (d) Imagine the app repeats the clustering analysis every week based on the modified data. Repeat this process 50 times. Plot the original data cloud, and the final data cloud you obtained side by side. Describe what has happened to your original population of diverse individuals after a year (50 weeks).

6. Presidential logorrhea (50 points)

In this problem, you will analyze the lengths of the State of the Union addresses.

- (a) The transcripts of all State of the Union addresses are linked on the canvas homework page or can be accessed at
`/project/cmsc25025/sou/speeches.pkl`.
Once you download them to your machine you can load them into python using the pickle package in python as follows:

```
import pickle
f=open('speeches.pkl','r')
speeches=pickle.load(f)
```

Write Python code that parses each SOU address, finding end-of-sentence markers. Don't worry about being too precise about sentence boundaries—as a first approximation, you could find words ending in a period. (But what about “Mr.”?)

- (b) For each year, compute the number of sentences in the address, and the mean sentence length in words for that year. Plot these data and two linear regressions, one plot for the number of sentences by year, another for the average sentence length by year. Note that the definition of “word” and “sentence” is imprecise. You can experiment with different parsing rules, and see if the results change qualitatively. Describe the trends that you see, and give some explanation for them. You should compute the linear regressions directly—for example, you may use the linear algebra routine `numpy.linalg.solve` but do not use a package that computes the regression.
- (c) Now, compute two regressions of the total number of words in a SOU versus year—one for the years 1790 to 1912, another for the years 1913 to the present. What trends do you see? Lookup the history of the State of the Union addresses (for example on Wikipedia) to explain the regressions.
- (d) Which President has the longest sentences on average? Which has the shortest sentences? Compute the median, 25% and 75% quantiles across all Presidents. What was the longest and shortest sentence ever spoken (or written) in a SOU?