

## Homework 9 Least squares and PCA

### Least squares

①

- Ⓐ Show that the best least-squares fit to a set of measurements  $y_1, \dots, y_m$  by a *horizontal line* (a constant function  $y = C$ ) is their average

$$C = \frac{y_1 + \dots + y_m}{m}.$$

Ⓑ

- Show that the slope of the line that passes through the origin in  $\mathbb{R}^2$  and comes closest in the least squares sense to passing through the points  $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$  is given by  $m = \sum_i x_i y_i / \sum_i x_i^2$ .

②

Find the best straight-line fit (least squares) to the measurements

$$\begin{array}{ll} b = 4 & \text{at } t = -2, \\ b = 1 & \text{at } t = 0, \end{array} \quad \begin{array}{ll} b = 3 & \text{at } t = -1, \\ b = 0 & \text{at } t = 2. \end{array}$$

Then find the projection of  $b = (4, 3, 1, 0)$  onto the column space of

$$A = \begin{bmatrix} 1 & -2 \\ 1 & -1 \\ 1 & 0 \\ 1 & 2 \end{bmatrix}.$$

③

- Find the best least-squares error *parabola* to the four points  $(0, 0), (1, 8), (3, 8), (4, 20)$ .  
(Your answer should minimize the summed squared error in the  $y$ -coordinates.) What is the  $R^2$  value of your fit?

### Economics

④

An economist hypothesizes that the change (in dollars) in the price of a loaf of bread is primarily a linear combination of the change in the price

An economist hypothesizes that the change (in dollars) in the price of a loaf of bread is primarily a linear combination of the change in the price of a bushel of wheat and the change in the minimum wage. That is, if  $B$  is the change in bread prices,  $W$  is the change in wheat prices, and  $M$  is the change in the minimum wage, then  $B = \alpha W + \beta M$ . Suppose that for three consecutive years the change in bread prices, wheat prices, and the minimum wage are as shown below.

	Year 1	Year 2	Year 3
$B$	+\$1	+\$1	+\$1
$W$	+\$1	+\$2	0\$
$M$	+\$1	0\$	-\$1

Use the theory of least squares to estimate the change in the price of bread in Year 4 if wheat prices and the minimum wage each fall by \$1.

⑤

Okuns "law", in economics, states that the annual change in gross domestic product (GDP) should relate to the annual change in the unemployment rate via an equation

$$\underset{\text{change in GDP}}{\Delta G} = k - c \cdot \underset{\text{change in unemployment rate}}{\Delta U}$$

I have uploaded two datasets, giving US GDP growth rates and unemployment rates. You can read them into Matlab using the `csvread()` function.

Find the least-squares best values for  $k$  and  $c$ .

Plot the data and your best-fit line.

Explain briefly what this means (one or two sentences).

(Note: The data gives the unemployment rates. You'll need to compute the changes in unemployment rates.)

Least squares and PCA in astronomy

# Least squares and PCA in astronomy

⑥

Henrietta Leavitt discovered thousands of "variable stars" (stars with varying brightness) in the Magellanic Clouds (dwarf galaxies orbiting our Milky Way galaxy). Among them were 25 stars with regular periods now known as "Cepheid variable stars." [https://en.wikipedia.org/wiki/Henrietta\\_Swan\\_Leavitt](https://en.wikipedia.org/wiki/Henrietta_Swan_Leavitt)

In 1912, she published a relationship between the period and the luminosity of Cepheids. (Since they were in the same cluster, the stars she studied were all roughly at the same distance from the Earth, so the measured luminosities were comparable.)

Cepheids became the first "standard candle" in astronomy, allowing astronomers to calculate distances to other galaxies. [https://en.wikipedia.org/wiki/Cosmic\\_distance\\_ladder](https://en.wikipedia.org/wiki/Cosmic_distance_ladder)

In fact, the Cepheids provided strong evidence that there **were** other galaxies. Using a telescope on Mt Wilson (here in LA!), Edwin Hubble discovered Cepheids, measured their brightness, and extrapolated their distances---far outside the Milky Way. <https://timesmachine.nytimes.com/timesmachine/1924/11/23/issue.html>

Here is her data:

TABLE I.  
PERIODS OF VARIABLE STARS IN THE SMALL MAGELLANIC CLOUD.

H.	Max.	Min.	Epoch.	Period.	Res. M.	Res. m.	H.	Max.	Min.	Epoch.	Period.	Res. M.	Res. m.
			d.	d.						d.	d.		
1505	14.8	16.1	0.02	1.25336	-0.6	-0.5	1400	14.1	14.8	4.0	6.650	+0.2	-0.3
1436	14.8	16.4	0.02	1.6637	-0.3	+0.1	1355	14.0	14.8	4.8	7.483	+0.2	-0.2
1446	14.8	16.4	1.38	1.7620	-0.3	+0.1	1374	13.9	15.2	6.0	8.397	+0.2	-0.3
1506	15.1	16.3	1.08	1.87502	+0.1	+0.1	818	13.6	14.7	4.0	10.336	0.0	0.0
1413	14.7	15.6	0.35	2.17352	-0.2	-0.5	1610	13.4	14.6	11.0	11.645	0.0	0.0
1460	14.4	15.7	0.00	2.913	-0.3	-0.1	1365	13.8	14.8	9.6	12.417	+0.4	+0.2
1422	14.7	15.9	0.6	3.501	+0.2	+0.2	1351	13.4	14.4	4.0	13.08	+0.1	-0.1
842	14.6	16.1	2.61	4.2897	+0.3	+0.6	827	13.4	14.3	11.6	13.47	+0.1	-0.2
1425	14.3	15.3	2.8	4.547	0.0	-0.1	822	13.0	14.6	13.0	16.75	-0.1	+0.3
1742	14.3	15.5	0.95	4.9866	+0.1	+0.2	823	12.2	14.1	2.9	31.94	-0.3	+0.4
1646	14.4	15.4	4.30	5.311	+0.3	+0.1	824	11.4	12.8	4.	65.8	-0.4	-0.2
1649	14.3	15.2	5.05	5.323	+0.2	-0.1	821	11.2	12.1	97.	127.0	-0.1	-0.4
1492	13.8	14.8	0.6	6.2926	-0.2	-0.4							

[http://adsbit.harvard.edu/cgi-bin/nph-article\\_query?1912HarCi.173...11&defaultprint=Y&filetype=pdf](http://adsbit.harvard.edu/cgi-bin/nph-article_query?1912HarCi.173...11&defaultprint=Y&filetype=pdf)

The New York Times  
SUNDAY, NOVEMBER 23, 1924

**FINDS SPIRAL NEBULAE ARE STELLAR SYSTEMS**  
**Dr. Hubble Confirms View That They Are 'Island Universes' Similar to Our Own.**

WASHINGTON, Nov. 22.—Confirmation of the view that the spiral nebulae, which appear in the heavens as whirling clouds, are in reality distant stellar systems, or "island universes," has been obtained by Dr. Edwin Hubble of the Carnegie Institution's Mount Wilson observatory, through investigations carried out with the observatory's powerful telescopes.

The number of spiral nebulae, the observatory officials have reported to the institution, is very great, amounting to hundreds of thousands, and their apparent sizes range from small objects, almost starlike in character, to the great nebulae in Andromeda, which extends across an angle some 2 degrees in the heavens, about six times the diameter of the full moon.

"The investigations of Dr. Hubble were made photographically with the 60-inch and 100-inch reflectors of the Mount Wilson observatory," the report said, "the extreme faintness of the stars under examination making necessary the use of these great telescopes. The resulting series of these instruments breaks up the outer portions of the nebulae into clouds of stars, which may be studied individually and compared with those in our own system."

From an investigation of the photographs thirty-six variable stars of the type referred to, known as Cepheid variables, were discovered in the two spirals, Andromeda and No. 33, of Messier's great catalogue of nebulae. The study of the periods of these stars and the application of the relationship between length of period and intrinsic brightness at once provided the means of determining the distances of these objects.

"The results are striking in their confirmation of the view that these spiral nebulae are distant stellar systems, and that we are observing them by light rays that have traveled as far as 1,000,000 light years. This means that light traveling at the rate of 186,000 miles a second has required a million years to reach us from these nebulae and that we are observing them by light rays that have traveled as far as 1,000,000 light years."

"With a knowledge of the distances of these nebulae we find for their diameters 4,500 light years for the Andromeda nebulae and 15,000 light years for Messier 33. These quantities, as well as the masses and densities of the systems, are quite comparable with the corresponding values for our local system of stars."

**FUNDS FOR SCHENCK HOUSE**  
**William C. Redfield Says It Was Built of Timbers of Old Ship.**

William C. Redfield, formerly Secretary of Commerce and now the President of the Netherlands-America Foundation, 17 East Forty-second Street, was one of the many who were interested in the news printed in yesterday's Times that an offer had been submitted to Murray Hulbert, President of the Board of Aldermen, to sell to the city for \$10,000 the old Schenck housestead at Hill Street, Brooklyn, which is believed to be the oldest house in New York City.

Mr. Redfield, in a letter to Mr. Hulbert yesterday, said that the Schenck house was built out of the timbers of an ancient ship. The old beams are visible and the bones of the old vessel still support the upper floors.

"I earnestly hope that funds may be made available, in order that this exceptional landmark of our city's history may be preserved," wrote Mr. Redfield. Mrs. Redfield is connected by marriage with the Schenck family.

Find the best-fitting lines for maximum luminosity as a function of log(period), and for minimum luminosity as a function of log(period). Plot your results. Also give the R<sup>2</sup> values.

⑦

In 1929, Edwin Hubble famously showed that the universe is expanding. Specifically, he showed a roughly linear relationship between the distances of other galaxies and their velocities away from us.

Here is the data he used:

NGC #	Distance (x10 <sup>6</sup> parsecs)	Radial velocity (km/sec)	Right ascension	Declination	Adjusted radial velocity
NA	0.032	170	NA	NA	170
NA	0.034	290	NA	NA	290
6822	0.214	-130	{19, 44, 57.8}	{-14, 48, 11}	60
598	0.263	-70	{1, 33, 51.}	{30, 39, 37}	15
221	0.275	-185	{0, 42, 41.9}	{40, 51, 57}	-30
224	0.275	-220	{0, 42, 44.3}	{41, 16, 9}	-65
5457	0.45	200	{14, 3, 12.5}	{54, 20, 53}	395
4736	0.5	290	{12, 50, 52.6}	{41, 7, 9}	405
5194	0.5	270	{13, 29, 52.4}	{47, 11, 41}	430
4449	0.63	200	{12, 28, 11.}	{44, 5, 33.4}	305
4214	0.8	300	{12, 15, 39.2}	{36, 19, 41}	370
3031	0.9	-30	{9, 55, 33.2}	{69, 3, 55}	90
3627	0.9	550	{11, 20, 15.1}	{12, 50, 22}	500

[https://en.wikipedia.org/wiki/Hubble%27s\\_law](https://en.wikipedia.org/wiki/Hubble%27s_law)

NGC #	Distance (x10 <sup>6</sup> parsecs)	Radial velocity (km/sec)	Right ascension	Declination	Adjusted radial velocity
NA	0.032	170	NA	NA	170
NA	0.034	290	NA	NA	290
6822	0.214	-130	{19, 44, 57.8}	{-14, 48, 11}	60
598	0.263	-70	{1, 33, 51.}	{30, 39, 37}	15
221	0.275	-185	{0, 42, 41.9}	{40, 51, 57}	-30
224	0.275	-220	{0, 42, 44.3}	{41, 16, 9}	-65
5457	0.45	200	{14, 3, 12.5}	{54, 20, 53}	395
4736	0.5	290	{12, 50, 52.6}	{41, 7, 9}	405
5194	0.5	270	{13, 29, 52.4}	{47, 11, 41}	430
4449	0.63	200	{12, 28, 11.}	{44, 5, 33.4}	305
4214	0.8	300	{12, 15, 39.2}	{36, 19, 41}	370
3031	0.9	-30	{9, 55, 33.2}	{69, 3, 55}	90
3627	0.9	650	{11, 20, 15.1}	{12, 59, 22}	580
4826	0.9	150	{12, 56, 43.9}	{21, 41, 0}	205
5236	0.9	500	{13, 37, 0.8}	{-29, 51, 59}	425
1068	1	920	{2, 42, 40.8}	{0, 0, 48}	830
5055	1.1	450	{13, 15, 49.3}	{42, 1, 47}	585
7331	1.1	500	{22, 37, 4.3}	{34, 24, 59}	740
4258	1.4	500	{12, 18, 57.5}	{47, 18, 14}	610
4151	1.7	960	{12, 10, 32.7}	{39, 24, 20}	1035
4382	2	500	{12, 25, 24.2}	{18, 11, 27}	515
4472	2	850	NA	NA	850
4486	2	800	{12, 30, 49.4}	{12, 23, 28}	800
4649	2	1090	{12, 43, 40.2}	{11, 33, 9}	1100

The second column gives the distance to each galaxy, and the last column gives the velocity away from us. (Hubble actually started with the data in column 3, but I have adjusted these velocities for the motion of our Sun.)

Ⓐ Run linear regression of distance versus velocity to find the best-fitting line. Make sure your line goes through (0,0)!!

Ⓑ Now run linear regression of velocity versus distance. Why does this give a different answer than part Ⓐ?

Ⓒ Now use PCA to find the best-fitting line.

Plot all three lines, and the data, on one labeled graph. Why is PCA more appropriate for analyzing this data than either linear regression?

Ⓓ Some of these data points are more precise than others. For example, they may have been collected by different telescopes.

In class, we saw how to get the  $k$ -dimensional subspace  $S$  that minimizes

$$\sum_{j=1}^m \|\vec{x}_j - P_S \vec{x}_j\|^2,$$

the sum of the squared distances from the data points to their projections on  $S$ . (The answer was to set  $S = \text{Span}\{k \text{ largest e-value e-vectors of } \sum x_j x_j^T\}$ .)

Extend this analysis to show how to get the  $k$ -dim<sup>l</sup> subspace  $S$  that minimizes

$$2 \cdot \|\vec{x}_1 - P_S \vec{x}_1\|^2 + \sum_{j=2}^m \|\vec{x}_j - P_S \vec{x}_j\|^2.$$

(This situation would arise if data point  $\vec{x}_1$  was more precise than the others.)

## Principle component analysis

⑧ The file "data.mat" contains roughly 15,000 data points in 32 dimensions. Use the singular-value decomposition (SVD) to project the data onto the best two-dimensional affine plane. Plot the projected data set.   
 (If you are using Matlab, the "scatter" function might be helpful for plotting.)

x-axis = principal component  
y-axis = second component