

to show that if Newton's method applied to $f(x)$ updates $x^{(i)}$ to $x^{(i+1)}$, then Newton's method applied to $g(z)$ will update $z^{(i)} = A^{-1}x^{(i)}$ to $z^{(i+1)} = A^{-1}x^{(i+1)}$.³

(b) [3 points] Is gradient descent invariant to linear reparameterizations? Justify your answer.

5. [35 points] Regression for denoising quasar spectra⁴

Introduction. In this problem, we will apply a supervised learning technique to estimate the light spectrum of *quasars*. Quasars are luminous distant galactic nuclei that are so bright, their light overwhelms that of stars in their galaxies. Understanding properties of the spectrum of light emitted by a quasar is useful for a number of tasks: first, a number of quasar properties can be estimated from the spectra, and second, properties of the regions of the universe through which the light passes can also be evaluated (for example, we can estimate the density of neutral and ionized particles in the universe, which helps cosmologists understand the evolution and fundamental laws governing its structure). The *light spectrum* is a curve that relates the light's intensity (formally, lumens per square meter), or *luminous flux*, to its wavelength. Figure 1 shows an example of a quasar light spectrum, where the wavelengths are measured in Angstroms (\AA), where $1\text{\AA} = 10^{-10}$ meters.

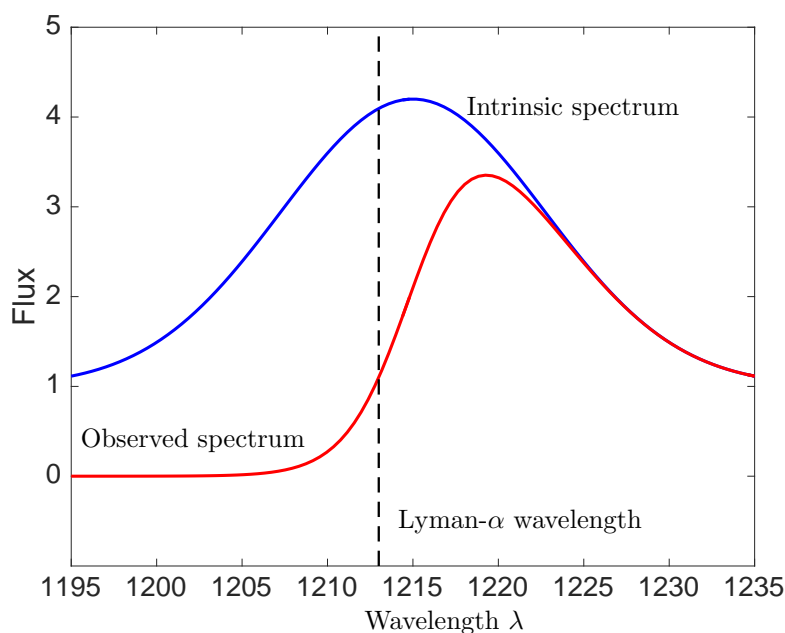


Figure 1: Light spectrum of a quasar. The blue line shows the intrinsic (i.e. original) flux spectrum emitted by the quasar. The red line denotes the observed spectrum here on Earth. To the left of the Lyman- α line, the observed flux is damped and the intrinsic (unabsorbed) flux continuum is not clearly recognizable (red line). To the right of the Lyman- α line, the observed flux approximates the intrinsic spectrum.

The Lyman- α wavelength is a wavelength beyond which intervening particles at most negligibly interfere with light emitted from the quasar. (Interference generally occurs when a photon is

³Note that for this problem, you must explicitly prove any matrix calculus identities that you wish to use that are not given in the lecture notes.

⁴Ciollaro, Mattia, et al. "Functional regression for quasar spectra." arXiv:1404.3168 (2014).

absorbed by a neutral hydrogen atom, which only occurs for certain wavelengths of light.) For wavelengths greater than this Lyman- α wavelength, the observed light spectrum f_{obs} can be modeled as a smooth spectrum f plus noise:

$$f_{\text{obs}}(\lambda) = f(\lambda) + \text{noise}(\lambda)$$

For wavelengths below the Lyman- α wavelength, a region of the spectrum known as the Lyman- α forest, intervening matter causes attenuation of the observed signal. As light emitted by the quasar travels through regions of the universe richer in neutral hydrogen, some of it is absorbed, which we model as

$$f_{\text{obs}}(\lambda) = \text{absorption}(\lambda) \cdot f(\lambda) + \text{noise}(\lambda)$$

Astrophysicists and cosmologists wish to understand the absorption function, which gives information about the Lyman- α forest, and hence the distribution of neutral hydrogen in otherwise unreachable regions of the universe. This gives clues toward the formation and evolution of the universe. Thus, it is our goal to estimate the spectrum f of an observed quasar.

Getting the data. We will be using data generated from the Hubble Space Telescope Faint Object Spectrograph (HST-FOS), Spectra of Active Galactic Nuclei and Quasars.⁵ We have provided two comma-separated data files located at:

- Training set: http://cs229.stanford.edu/ps/ps1/quasar_train.csv
- Test set: http://cs229.stanford.edu/ps/ps1/quasar_test.csv

Each file contains a single header row containing 450 numbers corresponding integral wavelengths in the interval $[1150, 1600]$ Å. The remaining lines contain relative flux measurements for each wavelength. Specifically, `quasar_train.csv` contains 200 examples and `quasar_test.csv` contains 50 examples. You may use the helper file `load_quasar_data.m` to load the data in Matlab: http://cs229.stanford.edu/ps/ps1/load_quasar_data.m

(a) [10 points] Locally weighted linear regression

Consider a linear regression problem in which we want to “weight” different training examples differently. Specifically, suppose we want to minimize

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m w^{(i)} \left(\theta^T x^{(i)} - y^{(i)} \right)^2$$

In class, we worked out what happens for the case where all the weights (the $w^{(i)}$'s) are the same. In this problem, we will generalize some of those ideas to the weighted setting.

- i. [2 points] Show that $J(\theta)$ can also be written

$$J(\theta) = (X\theta - \vec{y})^T W (X\theta - \vec{y})$$

for an appropriate diagonal matrix W , and where X and \vec{y} are as defined in class. State clearly what W is.

- ii. [4 points] If all the $w^{(i)}$'s equal 1, then we saw in class that the normal equation is

$$X^T X \theta = X^T \vec{y},$$

⁵<https://hea-www.harvard.edu/FOSAGN/>

and that the value of θ that minimizes $J(\theta)$ is given by $(X^T X)^{-1} X^T \vec{y}$. By finding the derivative $\nabla_{\theta} J(\theta)$ and setting that to zero, generalize the normal equation to this weighted setting, and give the new value of θ that minimizes $J(\theta)$ in closed form as a function of X , W and \vec{y} .

- iii. [4 points] Suppose we have a training set $\{(x^{(i)}, y^{(i)}); i = 1 \dots, m\}$ of m independent examples, but in which the $y^{(i)}$'s were observed with differing variances. Specifically, suppose that

$$p(y^{(i)} | x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi}\sigma^{(i)}} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2(\sigma^{(i)})^2}\right)$$

I.e., $y^{(i)}$ has mean $\theta^T x^{(i)}$ and variance $(\sigma^{(i)})^2$ (where the $\sigma^{(i)}$'s are fixed, known, constants). Show that finding the maximum likelihood estimate of θ reduces to solving a weighted linear regression problem. State clearly what the $w^{(i)}$'s are in terms of the $\sigma^{(i)}$'s.

- (b) [6 points] Visualizing the data

- i. [2 points] Use the normal equations to implement (unweighted) linear regression ($y = \theta^T x$) on the *first* training example (i.e. first non-header row). On one figure, plot both the raw data and the straight line resulting from your fit. State the optimal θ resulting from the linear regression. Remember the intercept term (your optimal θ should lie in \mathbb{R}^2).
- ii. [2 points] Implement locally weighted linear regression on the *first* training example. Use the normal equations you derived in part (a)(ii). On a different figure, plot both the raw data and the smooth curve resulting from your fit. When evaluating $h(\cdot)$ at a query point x , use weights

$$w^{(i)} = \exp\left(-\frac{(x - x^{(i)})^2}{2\tau^2}\right),$$

with bandwidth parameter $\tau = 5$.

- iii. [2 points] Repeat (b)(ii) four more times with $\tau = 1, 10, 100$ and 1000 . Plot the resulting curves. You can submit one plot with all four τ values or submit four separate plots. If you submit one plot, make sure all curves are visible. Additionally, in **2-3 sentences**, comment on what happens to the locally weighted linear regression line as τ varies.

- (c) [19 points] Predicting quasar spectra with functional regression

We now go a step beyond what we have covered explicitly in class, and we wish to predict an entire part of a spectrum—a curve—from noisy observed data. We begin by supposing that we observe a random sample of m absorption-free spectra, which is possible for quasars very close (in a sense relative to the size of the universe!) to Earth. For a given spectrum f , define f_{right} to be the spectrum to the right of the Lyman- α line. Let f_{left} be the spectrum within the Lyman- α forest region, that is, for lower wavelengths. To make the results cleaner, we define:

$$f(\lambda) = \begin{cases} f_{\text{left}}(\lambda) & \text{if } \lambda < 1200 \\ f_{\text{right}}(\lambda) & \text{if } \lambda \geq 1300 \end{cases}$$

We will learn a function r (for regression) that maps an observed f_{right} to an unobserved target f_{left} (note that f_{left} and f_{right} don't cover the entire spectrum). This is useful in practice because we observe f_{right} with *only* random noise: there is no systematic absorption, which we cannot observe directly, because hydrogen does not absorb photons with

higher wavelengths. By predicting f_{left} from a noisy version of f_{right} , we can estimate the unobservable spectrum of a quasar as well as the absorption function. Imaging systems collect data of the form

$$f_{\text{obs}}(\lambda) = \text{absorption}(\lambda) \cdot f(\lambda) + \text{noise}(\lambda)$$

for $\lambda \in \{\lambda_1, \dots, \lambda_n\}$, a *finite* number of points λ , because they must quantize the information. That is, even in the quasars-close-to-Earth training data, our observations of f_{left} and f_{right} consist of noisy evaluations of the true spectrum f at multiple wavelengths. In our case, we have $n = 450$ and $\lambda_1 = 1150, \dots, \lambda_n = 1599$.

We formulate the functional regression task as the goal of learning the function r mapping f_{right} to f_{left} :

$$r(f_{\text{right}})(\lambda) = \mathbb{E}(f_{\text{left}} \mid f_{\text{right}})(\lambda)$$

for λ in the Lyman- α forest.

- i. [1 points] First, we must smooth the data in the training dataset to make it more useful for prediction. For each $i = 1, \dots, m$, define $f^{(i)}(\lambda)$ to be the weighted linear regression estimate the i^{th} spectrum. Use your code from part (b)(ii) above to smooth all spectra in the training set using $\tau = 5$. Do the same for the test set. Apply smoothing to the entire spectrums (including both f_{left} and f_{right}) for both train and test. We will now operate on these smoothed spectra.
- ii. [14 points] Using your estimated regression functions $f^{(i)}$ for $i = 1, \dots, m$, we now wish to estimate the unobserved spectrum f_{left} of a quasar from its (noisy) observed spectrum f_{right} . To do so, we perform a weighted regression of the *locally weighted regressions*. In particular, given a new noisy spectrum observation:

$$f_{\text{obs}}(\lambda) = f(\lambda) + \text{noise}(\lambda) \quad \text{for } \lambda \in \{1300, \dots, 1599\}.$$

We define a metric d which takes as input, two spectra f_1 and f_2 , and outputs a scalar:

$$d(f_1, f_2) = \sum_i \left(f_1(\lambda_i) - f_2(\lambda_i) \right)^2.$$

The metric d computes squared distance between the new datapoint and previous datapoints. If f_1 and f_2 are right spectra, then we take the preceding sum only over $\lambda \in \{1300, \dots, 1599\}$, rather than the entire spectrum.

Based on this distance function, we may define the nonparametric *functional* regression estimator, which is a locally weighted sum of *functions* f_{left} from the training data (this is like locally weighted linear regression, except that instead of predicting $y \in \mathbb{R}$ we predict a function f_{left}). Specifically, let f_{right} denote the right side of a spectrum, which we have smoothed using locally weighted linear regression (as you were told to do in the previous part of the problem). We wish to estimate the associated *left* spectrum f_{left} . Define the function $\text{ker}(t) = \max\{1 - t, 0\}$ and let $\text{neighb}_k(f_{\text{right}})$ denote the k indices $i \in \{1, 2, \dots, m\}$ of the training set that are closest to f_{right} , that is

$$d(f_{\text{right}}^{(i)}, f_{\text{right}}) < d(f_{\text{right}}^{(j)}, f_{\text{right}}) \quad \text{for all } i \in \text{neighb}_k(f_{\text{right}}), j \notin \text{neighb}_k(f_{\text{right}})$$

and $\text{neighb}_k(f_{\text{right}})$ contains exactly k indices. In addition, let

$$h := \max_{i \in \{1, \dots, m\}} d(f_{\text{right}}^{(i)}, f_{\text{right}}).$$

Then define the estimated function $\widehat{f_{\text{left}}} : \mathbb{R} \rightarrow \mathbb{R}$ by

$$\widehat{f_{\text{left}}}(\lambda) = \frac{\sum_{i \in \text{neighb}_k(f_{\text{right}})} \ker(d(f_{\text{right}}^{(i)}, f_{\text{right}})/h) f_{\text{left}}^{(i)}(\lambda)}{\sum_{i \in \text{neighb}_k(f_{\text{right}})} \ker(d(f_{\text{right}}^{(i)}, f_{\text{right}})/h)}. \quad (1)$$

Include f_{right} from training in its own neighborhood. Recall that $f_{\text{right}}^{(i)}$ is the *smoothed* (weighted linear regression) estimate of the i th training spectrum.

Construct the functional regression estimate (1) for each spectrum in the entire training set using $k = 3$ nearest neighbors: for each $j = 1, \dots, m$, construct the estimator $\widehat{f_{\text{left}}}$ from (1) using $f_{\text{right}} = f_{\text{right}}^{(j)}$. Then compute the error $d(f_{\text{left}}^{(j)}, \widehat{f_{\text{left}}})$ between the true spectrum $f_{\text{left}}^{(j)}$ and your estimated spectrum $\widehat{f_{\text{left}}}$ for each j , and return the average over the training data. What is your average training error?

- iii. [4 points] Perform functional regression on the test set using the same procedure as in the previous subquestion. Note: You must use neighbors $f_{\text{right}}^{(i)}, f_{\text{left}}^{(i)}$ from the **training set**, and f_{right} in the test set to predict the corresponding f_{left} in the test set. What is your average test error? For test examples 1 and 6, include a plot with both the entire smooth spectrum and the fitted curve $\widehat{f_{\text{left}}}$ curve on the same graph. You should submit two plots: one for test example 1 and one for test example 6.

Reminder: Please include in your submission a copy of your code and figures for the programming questions.