

Machine Learning and Computational Statistics, Spring 2016

Homework 1: Ridge Regression and SGD

Due: Friday, February 5, 2015, at 6pm (Submit via NYU Classes)

Instructions: Your answers to the questions below, including plots and mathematical work, should be submitted as a single PDF file. You may include your code inline or submit it as a separate file. You may either scan hand-written work or, preferably, write your answers using software that typesets mathematics (e.g. L^AT_EX, L^AT_EX, or MathJax via iPython).

1 Introduction

In this homework you will implement ridge regression using gradient descent and stochastic gradient descent. We've provided a lot of support Python code to get you started on the right track. References below to particular functions that you should modify are referring to the support code, which you can download from the website. If you have time after completing the assignment, you might pursue some of the following:

- Study up on numpy's "broadcasting" to see if you can simplify and/or speed up your code: <http://docs.scipy.org/doc/numpy/user/basics.broadcasting.html>
- Think about how you could make the code more modular so that you could easily try different loss functions and step size methods.
- Experiment with more sophisticated approaches to setting the step sizes for SGD (e.g. try out the recommendations in "Bottou's SGD Tricks" on the website)
- Investigate what happens to the convergence rate when you intentionally make the feature values have vastly different orders of magnitude. Try a dataset (could be artificial) where $\mathcal{X} \subset \mathbf{R}^2$ so that you can plot the convergence path of GD and SGD. Take a look at <http://imgur.com/a/Hqolp> for inspiration.
- Instead of taking 1 data point at a time, as in SGD, try minibatch gradient descent, where you use multiple points at a time to get your step direction. How does this effect convergence speed? Are you getting computational speedup as well by using vectorized code?
- Advanced: What kind of loss function will give us "quantile regression"?

Include any investigations you do in your submission, and we may award optional credit.

I encourage you to develop the habit of asking "what if?" questions and then seeking the answers. I guarantee this will give you a much deeper understanding of the material (and likely better performance on the exam and job interviews, if that's your focus). You're also encouraged to post your interesting questions on Piazza under "questions."

2 Linear Regression

2.1 Feature Normalization

When feature values differ greatly, we can get much slower rates of convergence of gradient-based algorithms. Furthermore, when we start using regularization, features with larger values can have a much greater effect on the final output for the same regularization cost – in effect, features with larger values become more important once we start regularizing. One common approach to feature normalization is to linearly transform (i.e. shift and rescale) each feature so that all feature values in the training set are in $[0, 1]$. Each feature gets its own transformation. We then apply the same transformations to each feature on the test set. It's important that the transformation is “learned” on the training set, and then applied to the test set. It is possible that some transformed test set values will lie outside the $[0, 1]$ interval.

Modify function `feature_normalization` to normalize all the features to $[0, 1]$. (Can you use numpy's “broadcasting” here?)

2.2 Gradient Descent Setup

In linear regression, we consider the hypothesis space of linear functions $h_\theta : \mathbf{R}^d \rightarrow \mathbf{R}$, where

$$h_\theta(x) = \theta^T x,$$

for $\theta, x \in \mathbf{R}^d$, and we choose θ that minimizes the following “square loss” objective function:

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x_i) - y_i)^2,$$

where $(x_1, y_1), \dots, (x_m, y_m) \in \mathbf{R}^d \times \mathbf{R}$ is our training data.

While this formulation of linear regression is very convenient, it's more standard to use a hypothesis space of “affine” functions:

$$h_{\theta,b}(x) = \theta^T x + b,$$

which allows a “bias” or nonzero intercept term. The standard way to achieve this, while still maintaining the convenience of the first representation, is to add an extra dimension to x that is always a fixed value, such as 1. You should convince yourself that this is equivalent. We'll assume this representation, and thus we'll actually take $\theta, x \in \mathbf{R}^{d+1}$.

1. Let $X \in \mathbf{R}^{m \times d+1}$ be the “design matrix”, where the i 'th row of X is x_i . Let $y = (y_1, \dots, y_m)^T \in \mathbf{R}^{m \times 1}$ be the “response”. Write the objective function $J(\theta)$ as a matrix/vector expression, without using an explicit summation sign.

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x_i) - y_i)^2,$$

$$J(\theta) = \frac{1}{2m} (X\theta - y)^T (X\theta - y),$$

2. Write down an expression for the gradient of J .

$$\frac{\partial J(\theta)}{\partial \theta_k} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i) x_i^k,$$

In matrix form

$$\nabla_{\theta} J(\theta) = \frac{1}{m} (X\theta - y)^T X,$$

3. In our search for a θ that minimizes J , suppose we take a step from θ to $\theta + \eta\Delta$, where $\Delta \in \mathbf{R}^{d+1}$ is a unit vector giving the direction of the step, and $\eta \in \mathbf{R}$ is the length of the step. Use the gradient to write down an approximate expression for $J(\theta + \eta\Delta) - J(\theta)$. [This approximation is called a “linear” or “first-order” approximation.]

The equation of a line with point (x_1, y_1) and slope m

$$(y - y_1) = m(x - x_1)$$

Here slope is given by $\frac{\partial J(\theta)}{\partial \theta}$ Substituting values in line question we get

$$J(\theta + \eta\Delta) - J(\theta) = \frac{\partial J(\theta)}{\partial \theta_k} (\eta\Delta)$$

The value of $\frac{\partial J(\theta)}{\partial \theta_k}$ is derived in part 2 above.

4. Write down the expression for updating θ in the gradient descent algorithm. Let η be the step size.

$$\theta_k = \theta_k - \eta \frac{\partial J(\theta)}{\partial \theta_k}$$

The value of $\frac{\partial J(\theta)}{\partial \theta_k}$ is derived in part 2 above.

5. Modify the function `compute_square_loss`, to compute $J(\theta)$ for a given θ .
6. Create a small dataset for which you can compute $J(\theta)$ by hand, and verify that your `compute_square_loss` function returns the correct value.

Consider the following:

$$X = ([1, 2, 3, 1], [0, 1, 7, 1]) \quad \theta = ([3, 4, 1, 2]) \quad \text{and} \quad y = ([20, 25])$$

We get $J(\theta) = 40$.

The same is shown in the code.

7. Modify the function `compute_square_loss_gradient`, to compute $\nabla_{\theta} J(\theta)$.
8. Using your small dataset, verify that your `compute_square_loss_gradient` function returns the correct value.

Considering the same values for X , θ and y as in 2.2.6, we get

$$\nabla_{\theta} J(\theta) = [-2, -10, -48, -8]$$

The same is verified in the code.

2.3 Gradient Checker

For many optimization problems, coding up the gradient correctly can be tricky. Luckily, there is a nice way to numerically check the gradient calculation. If $J : \mathbf{R}^d \rightarrow \mathbf{R}$ is differentiable, then for any direction vector $\Delta \in \mathbf{R}^d$, the directional derivative of J at θ in the direction Δ is given by¹:

$$\lim_{\epsilon \rightarrow 0} \frac{J(\theta + \epsilon \Delta) - J(\theta - \epsilon \Delta)}{2\epsilon}$$

We can approximate this directional derivative by choosing a small value of $\epsilon > 0$ and evaluating the quotient above. We can get an approximation to the gradient by approximating the directional derivatives in each coordinate direction and putting them together into a vector. In other words, take $\Delta = (1, 0, 0, \dots, 0)$ to get the first component of the gradient. Then take $\Delta = (0, 1, 0, \dots, 0)$ to get the second component. And so on. See http://ufldl.stanford.edu/wiki/index.php/Gradient_checking_and_advanced_optimization for details.

1. Complete the function `grad_checker` according to the documentation given.
2. (Optional) Write a generic version of `grad_checker` that will work for any objective function. It should take as parameters a function that computes the objective function and a function that computes the gradient of the objective function.

2.4 Batch Gradient Descent

At the end of the skeleton code, the data is loaded, split into a training and test set, and normalized. We'll now finish the job of running regression on the training set. Later on we'll plot the results together with SGD results.

1. Complete `batch_gradient_descent`.
2. Now let's experiment with the step size. Note that if the step size is too large, gradient descent may not converge². Starting with a step-size of 0.1, try various different fixed step sizes to see which converges most quickly and/or which diverge. As a minimum, try step sizes 0.5, 0.1, .05, and .01. Plot the value of the objective function as a function of the number of steps for each step size. Briefly summarize your findings.

With step size = 0.5, the initial loss was close to 441.4 and it went higher and higher, which means the step size taken is very large. The next θ value we taking is going beyond the minimum point. This diverges.

With step size = 0.05, the graph converges but quickly.

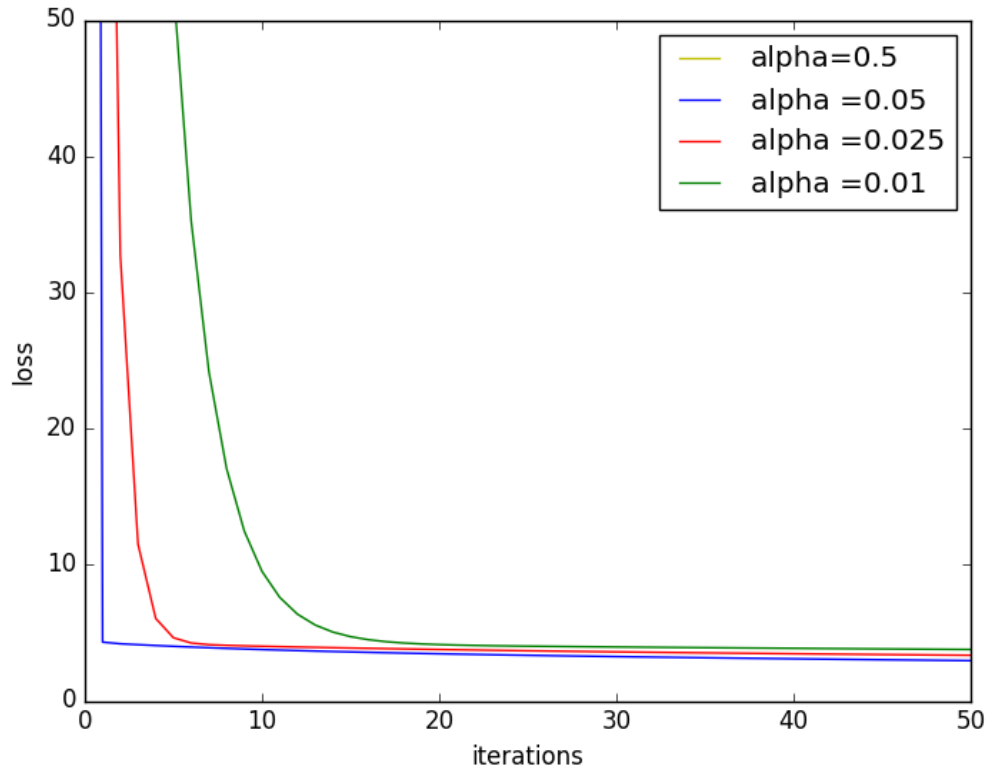
With step size = 0.025, the convergence is slower compared to step size = 0.05. This implies lower the step size, slower is the speed of convergence.

With step size = 0.01, the graph converges very slowly. This can be seen as green line that takes more iterations to converge.

¹Of course, it is also given by the more standard definition of directional derivative, $\lim_{\epsilon \downarrow 0} \frac{1}{\epsilon} [J(\theta + \epsilon \Delta) - J(\theta)]$. The form given gives a better approximation, but it requires differentiability at θ . It won't give the right result for a nondifferentiable function, such as $J(\theta) = |\theta|$ at $\theta = 0$ for $\Delta = 1$.

²For the mathematically inclined, there is a theorem that if the objective function is convex, differentiable, and Lipschitz continuous with constant $L > 0$, then gradient descent converges for fixed step sizes smaller than $1/L$. See https://www.cs.cmu.edu/~ggordon/10725-F12/scribes/10725_Lecture5.pdf, Theorem 5.1.

All these are depicted in figure below.



3. (Optional, but recommended) Implement backtracking line search (google it), and never have to worry choosing your step size again. How does it compare to the best fixed step-size you found in terms of number of steps? In terms of time? How does the extra time to run backtracking line search at each step compare to the time it takes to compute the gradient? (You can also compare the operation counts.)

2.5 Ridge Regression (i.e. Linear Regression with L_2 regularization)

When we have large number of features compared to instances, regularization can help control overfitting. Ridge regression is linear regression with L_2 regularization. The objective function is

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2 + \lambda \theta^T \theta,$$

where λ is the regularization parameter, which controls the degree of regularization. Note that the bias term is being regularized as well. We will address that below.

1. Compute the gradient of $J(\theta)$ and write down the expression for updating θ in the gradient descent algorithm.

$J(\theta)$ can be written as

$$J(\theta) = \frac{1}{2m}(X\theta - y)^T(X\theta - y) + \lambda\theta^T\theta,$$

Differentiating

$$\nabla_{\theta}J(\theta) = \frac{1}{m}(X\theta - y)^TX + 2\lambda\theta,$$

The update rule for θ can be written as

$$\theta_k = \theta_k - \eta \frac{\partial J(\theta)}{\partial \theta_k}$$

In matrix form

$$\theta = \theta - \eta \nabla_{\theta}J(\theta)$$

2. Implement `compute_regularized_square_loss_gradient`.
3. Implement `regularized_grad_descent`.
4. For regression problems, we may prefer to leave the bias term unregularized. One approach is to rewrite $J(\theta)$ and re-compute $\nabla_{\theta}J(\theta)$ in a way that separates out the bias from the other parameter. Another approach that can achieve approximately the same thing is to use a very large number B , rather than 1, for the extra bias dimension. Explain why making B large decreases the effective regularization on the bias term, and how we can make that regularization as weak as we like (though not zero).

The idea of regularization is to avoid the problem of over fitting. To avoid over fitting we need to make the curve smoother. This can be achieved by reducing the magnitude of θ excluding the bias term.

However, if we have to include bias term in regularization, it means we will be reducing its magnitude by regularization. The only way to keep its effect low is by increasing the initial value.

5. Choosing a reasonable step-size or using backtracking line search, find the θ_{λ}^* that minimizes $J(\theta)$ for a range of λ . You should plot the training loss and the validation loss (just the square loss part, without the regularization) as a function of λ . Your goal is to find λ that gives the minimum validation loss. It's hard to predict what λ that will be, so you should start your search very broadly, looking over several orders of magnitude. For example, $\lambda \in \{10^{-7}, 10^{-5}, 10^{-3}, 10^{-1}, 1, 10, 100\}$. Once you find a range that works better, keep zooming in. You may want to have $\log(\lambda)$ on the x -axis rather than λ .

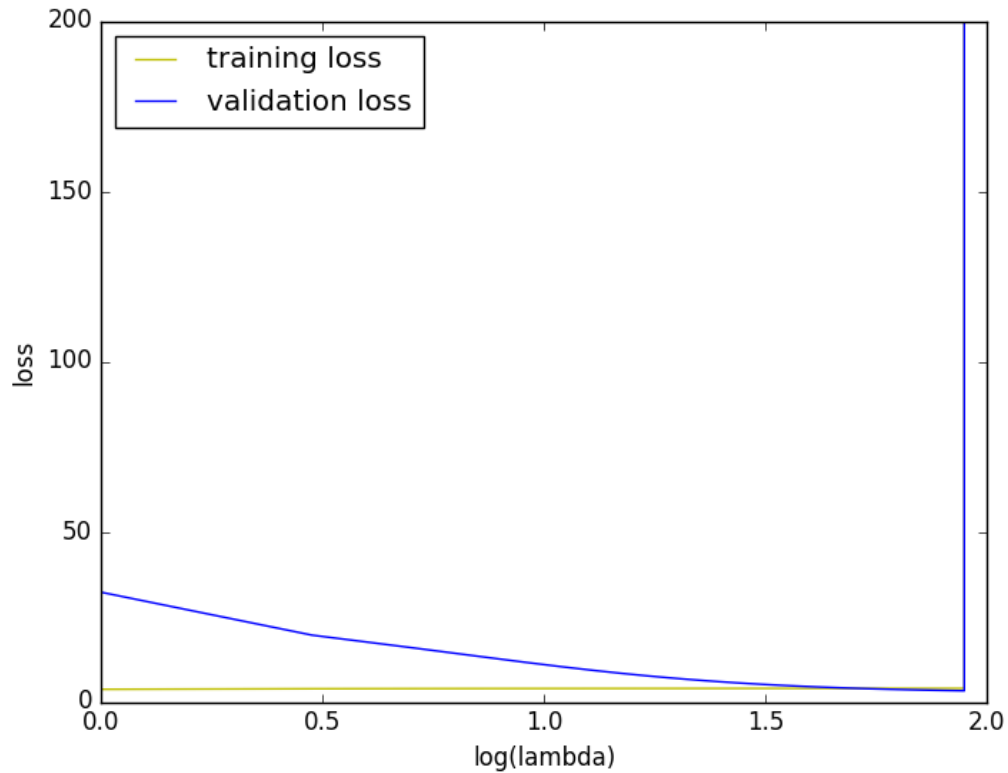
Selected a range of values from 10^{-10} to 100 and found that the minimum square loss is when the $\lambda = 10$. All this was calculated with step size = 0.01.

This is shown in graph below.

X-axis is $\log(\lambda)$ and Y axis is the loss.

Loss at $\lambda = 10$ is 10.96. This reduces even further till around $\lambda = 70$ but to make calculations easier and since the jump is pretty high after 70, to make things easier and be safe, let's consider $\lambda = 10$.

Things to be noted here is that validation loss decreases with λ value and then increases.



6. (Optional) Once you have found a good value for λ , repeat the fits with different values for B , and plot the results. For this dataset, does regularizing the bias help, hurt, or make no significant difference?
7. (Optional) Estimate the average time it takes on your computer to compute a single gradient step.
8. What θ would you select for deployment and why?

We found that $\lambda = 10$ and step size = 0.01 produced less validation loss. Hence corresponding θ obtained by calling `regularized_grad_descent` is used. We used validation loss to determine θ because our function is not trained using the test data and hence best results with test data is used for determining best θ

Value of θ used is

[-9.92635627e-03 6.67543885e-05 6.40706144e-03 8.43709881e-03 -1.42228569e-03 -1.34379919e-04 1.32553192e-03 1.32553192e-03 7.26646711e-03 1.27152013e-02 1.42641663e-02 6.47452075e-03 -2.25259582e-03 -7.04382209e-03 1.37839141e-02 1.64218885e-02 1.33508992e-02 2.75471484e-03 6.56536406e-04 6.56536406e-04 6.56536406e-04 2.69293752e-03 2.69293752e-03 2.69293752e-03 3.27840986e-03 3.27840986e-03 3.27840986e-03 3.54451035e-03 3.54451035e-03 3.54451035e-03 3.69052930e-03 3.69052930e-03 3.69052930e-03 5.25843606e-03 5.25843606e-03 5.25843606e-03 5.78340748e-03 5.78340748e-03 5.78340748e-03 5.05211594e-03 5.05211594e-03 5.05211594e-03 4.71122269e-03 4.71122269e-03 4.71122269e-03 4.52134806e-03 4.52134806e-03 4.52134806e-03 -1.28522936e-02]

2.6 Stochastic Gradient Descent

When the training data set is very large, evaluating the gradient of the loss function can take a long time, since it requires looking at each training example to take a single gradient step. In this case, stochastic gradient descent (SGD) can be very effective. In SGD, the gradient of the risk is approximated by a gradient at a single example. The approximation is poor, but it is unbiased. The algorithm sweeps through the whole training set one by one, and performs an update for each training example individually. One pass through the data is called an *epoch*. Note that each epoch of SGD touches as much data as a single step of batch gradient descent. Before we begin cycling through the training examples, it is important to shuffle them into a random order. You can use the same ordering for each epoch, though optionally you could investigate whether reshuffling after each epoch speeds up convergence.

1. Write down the update rule for θ in SGD.

$J(\theta)$ for SGD can be written as

$$J(\theta) = \frac{1}{2}(x_k\theta - y_k)^T(x_k\theta - y_k) + \lambda\theta^T\theta,$$

Where x_k and y_k is one the sample choosen randomly from the all the input samples.

Differentiating

$$\nabla_{\theta}J(\theta) = (x_k\theta - y_k)x_k + 2\lambda\theta,$$

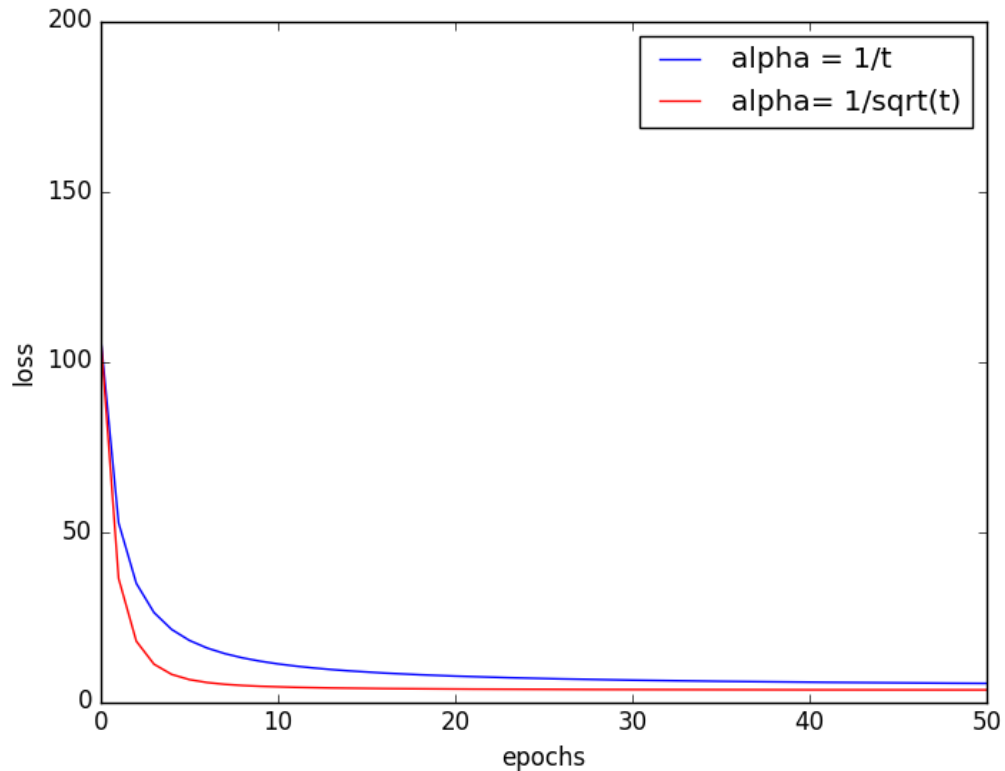
The update rule for θ can be written as

$$\theta = \theta - \eta\nabla_{\theta}J(\theta)$$

2. Implement `stochastic_grad_descent`.
3. Use SGD to find θ_{λ}^* that minimizes the ridge regression objective for the λ and B that you selected in the previous problem. (If you could not solve the previous problem, choose $\lambda = 10^{-2}$ and $B = 1$). Try a few fixed step sizes (at least try $\eta_t \in \{0.05, .005\}$). Note that SGD may not converge with fixed step size. Simply note your results. Next try step sizes that decrease with the step number according to the following schedules: $\eta_t = \frac{1}{t}$ and $\eta_t = \frac{1}{\sqrt{t}}$. For each step size rule, plot the value of the objective function (or the log of the objective function if that is more clear) as a function of epoch (or step number) for each of the approaches to step

size. How do the results compare? (Note: In this case we are investigating the convergence rate of the optimization algorithm, thus we're interested in the value of the objective function, which includes the regularization term.)

Graph for loss with respect to epochs is shown below for step size = $1/t$ and step size = $1/\sqrt{t}$. Note that, for step size = $1/t$, the graph converges slowly compared to step size = $1/\sqrt{t}$. This is because with step size = $1/t$, the step size reduces quickly. And lower value of step size implies slower rate of convergence. Also to be noted is that with SGD the graph converges much faster compared to GD.

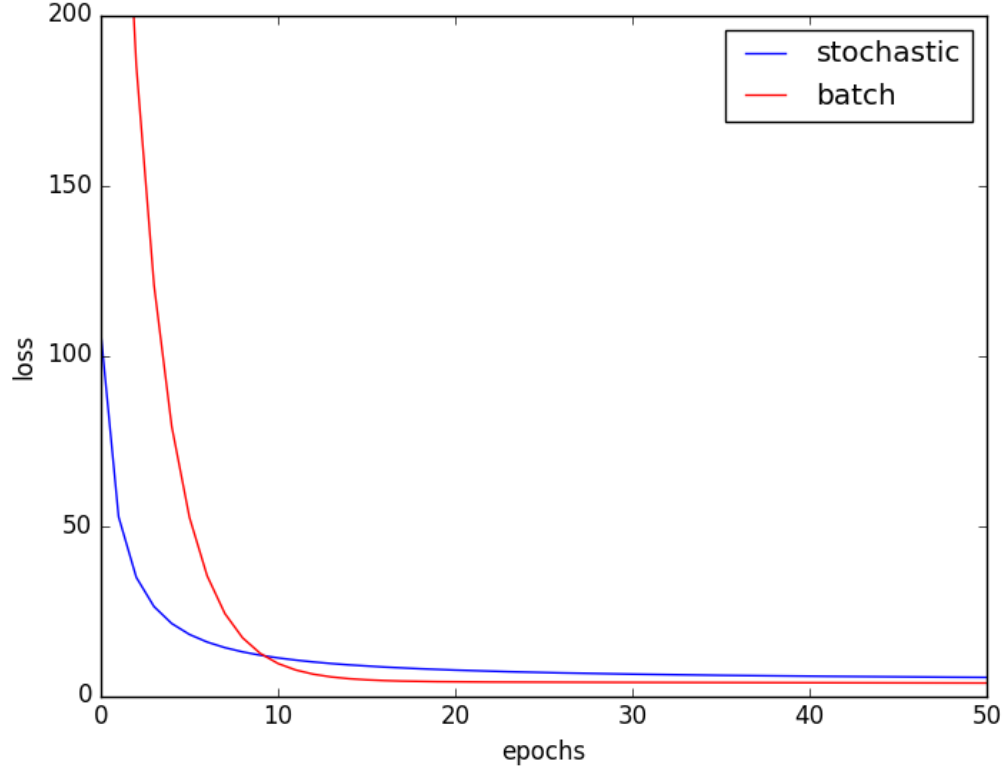


4. (Optional) Try a stepsize rule of the form $\eta_t = \frac{\eta_0}{1+\eta_0\lambda t}$, where λ is your regularization constant, and η_0 a constant you can choose. How do the results compare?
5. Estimate the amount of time it takes on your computer for a single epoch of SGD.
The amount of time it took to run one epoch is close to 0.001 second.
6. Comparing SGD and gradient descent, if your goal is to minimize the total number of epochs (for SGD) or steps (for batch gradient descent), which would you choose? If your goal were to minimize the total time, which would you choose?

Note that even though the step size in stochastic gradient descent reduces (initial for both is 0.01), stochastic converges faster.

This means in order to minimize total steps for convergence then stochastic gradient descent is better.

However, in order to minimize total time gradient descent is better because SGD computes one at a time. If we have to compute for all input values in an epoch then it takes longer time.



3 Risk Minimization

Recall that the definition of the **expected loss** or “**risk**” of a decision function $f : \mathcal{X} \rightarrow \mathcal{A}$ is

$$R(f) = \mathbb{E}\ell(f(x), y),$$

where $(x, y) \sim P_{\mathcal{X} \times \mathcal{Y}}$, and the **Bayes decision function** $f^* : \mathcal{X} \rightarrow \mathcal{A}$ is a function that achieves the *minimal risk* among all possible functions:

$$R(f^*) = \inf_f R(f).$$

Here we consider the regression setting, in which $\mathcal{A} = \mathcal{Y} = \mathbf{R}$.

1. Show that for the square loss $\ell(\hat{y}, y) = \frac{1}{2} (y - \hat{y})^2$, the Bayes decision function is a $f^*(x) = \mathbb{E}[Y | X = x]$. [Hint: Consider constructing $f^*(x)$, one x at a time.]

Square loss function can be written as

$$\ell(f(x), y) = \frac{1}{2} (f(x) - y)^2,$$

The risk of a square loss function will be

$$R(f) = \mathbb{E} \frac{1}{2} (f(x) - y)^2,$$

$$R(f) = \mathbb{E} \frac{1}{2} (f(x)^2 + y^2 - 2f(x)y)$$

Risk function is minimum for bayes decision function. Hence differentiating with respect to $f(x)$ and equating to 0 should give us $f(x)$ = bayes decision function.

Differentiating w.r.t $f(x)$, we get

$$\frac{\partial R(f)}{\partial f(x)} = \mathbb{E}(f(x) - y) = 0$$

Considering one input at a time, we can see that the bayes function is given by

$$f^*(x) = \mathbb{E}[Y | X = x]$$

2. (Optional) Show that for the absolute loss $\ell(\hat{y}, y) = |y - \hat{y}|$, the Bayes decision function is a $f^*(x) = \text{median}[Y | X = x]$. [Hint: Again, consider one x at time. For some approaches, it may help to use the following characterization of a median: m is a median of the distribution for random variable Y if $P(Y \geq m) \geq \frac{1}{2}$ and $P(Y \leq m) \geq \frac{1}{2}$.] Note: This loss function leads to “median regression”, There are other loss functions that lead to “quantile regression” for any chosen quantile.