

CS 475 Machine Learning (Spring 2017): Assignment 1

Due on February 27th, 2017 at 3:00PM

Raman Arora

Instructions: Please read these instructions carefully and follow them precisely. Feel free to ask the instructor if anything is unclear!

1. Please submit your solutions electronically via [Gradescope](#).
2. Please submit a PDF file for the written component of your solution including derivations, explanations, etc. You can create this PDF in any way you want: typeset the solution in LATEX (recommended), type it in Word or a similar program and convert/export to PDF, or even hand write the solution (legibly!) and scan it to PDF. We recommend that you restrict your solutions to the space allocated for each problem; you may need to adjust the white space by tweaking the argument to `\vspace{xpt}` command. Please name this document `<firstname-lastname>-sol1.pdf`.
3. Submit the empirical component of the solution (Python code and the documentation of the experiments you are asked to run, including figures) in a Jupyter notebook file.
4. **Late submissions:** You have a total of 72 late hours for the entire semester that you may use as you deem fit. After you have used up your quota, there will be a penalty of 50% of your grade on a late homework if submitted within 48 hours of the deadline and a penalty of 100% of your grade on the homework for submissions that are later than 48 hours past the deadline.
5. **What is the required level of detail?** When asked to derive something, please clearly state the assumptions, if any, and strive for balance: justify any non-obvious steps, but try to avoid superfluous explanations. When asked to plot something, please include the figure as well as the code used to plot it (and clearly explain in the README what the relevant files are). If multiple entities appear on a plot, make sure that they are clearly distinguishable (by color or style of lines and markers). When asked to provide a brief explanation or description, try to make your answers concise, but do not omit anything you believe is important. When submitting code, please make sure it's reasonably documented, and describe succinctly in the written component of the solution what is done in each `py`-file.

Name: _____

I. Regression

In this set of problems we will look at the regression problem and the maximum likelihood (ML) approach, with the goal to understand a bit better some of their properties.

We will start with simple linear regression, defined by

$$\hat{y} = f(\mathbf{x}; \mathbf{w}) = \mathbf{w} \cdot \mathbf{x} = \sum_{j=0}^d w_j x_j \quad (1)$$

It is assumed in (1) that we have augmented the inputs $\mathbf{x} \in \mathbb{R}^d$ by adding 1 as the “zeroth” dimension, $x_0 \equiv 1$. This assumption is valid for the remainder of this problem set, unless stated otherwise.

Assume that we are given n training examples $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)$, where each input data point \mathbf{x}_i has d real-valued features. The goal of regression is to learn to predict y from \mathbf{x} . We refer to $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{n \times d}$ as the design matrix and $\mathbf{y} = [y_1, \dots, y_n]^\top$ as the response vector.

Let $\hat{\mathbf{w}}$ be the linear regression parameters estimated using the least squares procedure from data, and \mathbf{w}^* be the best possible linear regression parameters for the underlying $p(\mathbf{x}, y)$.

It was shown in the lectures that prediction errors made by $\hat{\mathbf{w}}$ are uncorrelated with the training $\{\mathbf{x}_i\}$, and a more general statement was made but not proven: that the prediction errors are in fact uncorrelated with values of any linear function of \mathbf{x} computed on the training inputs.

Later, we made use of a similar fact regarding the prediction errors made by \mathbf{w}_* . Here we will prove this more general fact. First, to make things precise, let us recall the definition of correlation between two (continuous) random variables U and V . Let $p_{u,v}$ be their joint probability density, and p_u and p_v the corresponding marginal densities. Then, the correlation (coefficient) between U and V is defined as

$$\rho(U, V) \equiv \frac{\mathbb{E}_{p_{u,v}}[(U - \mathbb{E}_{p_u}[U])(V - \mathbb{E}_{p_v}[V])]}{\sqrt{\text{var}(U) \text{var}(V)}} \quad (2)$$

where $\text{var}(U)$ and $\text{var}(V)$ are the variances of U and V under their respective marginal distributions. Intuitively, correlation measures how well one variable is linearly predictable from the other. We will now prove a fact from which it follows that prediction errors of \mathbf{w}_* are uncorrelated with any linear function of \mathbf{x} :

Problem 1 [15 points] Show rigorously that for any $a \in \mathbb{R}^{d+1}$,

$$\mathbb{E}_{p(\mathbf{x}, y)} \left[(y - \mathbf{w}_*^\top \mathbf{x}) \mathbf{a}^\top \mathbf{x} \right] = 0 \quad (3)$$

Advice: You want to write explicitly what the definition of \mathbf{w}_* as the best linear regression under $p(\mathbf{x}, y)$ implies, in terms of minimizing the expected loss (which, to remind you, is an integral w.r.t. \mathbf{x} and y) similarly to how we treated the empirical loss in the case of $\hat{\mathbf{w}}$.

Let us begin by defining true risk as a function of the parameter w such that $R(w) = E_{(x_o, y_o) \sim p(x, y)}[\ell(f(x_o, w), y_o)]$ where $\ell(f(x_o, w), y_o)$ is the loss function. We can expand this expectation in terms of the marginal densities of x_o and y_o such that, $R(w) = E_{x_o \sim p(x)}[E_{y_o \sim p(y|x)}[\ell(f(x_o, w), y_o)|x_o]]$. By the definition of expectation, we can rewrite this as $\int_{x_o} E_{y_o \sim p(y|x)}[\ell(f(x_o, w), y_o)|x_o] p(x_o) dx_o$. For the case of linear regression, assume a least squares loss function, $\ell(f(x, w), y) = \sum_{i=1}^N (y - w \cdot x_i)^2$. Now, to minimize the true risk, we must pick the value of w that minimizes the inner conditional expectation, which is equivalent to picking the optimal value of w , w^* , such that $\ell(f(x, w), y)$ is minimized. Now, we get that $R(w) = \int_{x_o} E_{y_o \sim p(y|x)}[\sum_{i=1}^N (y - w \cdot x_i)^2 | x_o] p(x_o) dx_o$. However, can give this in matrix product form, $R(w) = \int_{x_o} E_{y_o \sim p(y|x)}[(y - w \cdot x)^2 | x_o] p(x_o) dx_o$. Now, we can differentiate $R(w)$ with respect to w : $\frac{\partial R(w)}{\partial w} = \frac{\partial}{\partial w} \int_{x_o} E_{y_o \sim p(y|x)}[(y - w^T x)^2 | x_o] p(x_o) dx_o$. Since we showed that we must only minimize the inner conditional to minimize $R(w)$, we can push the differentiation inside the integral.

$$\frac{\partial}{\partial w} E_{y \sim p(y|x)}[(y - w_*^T x)^2 | x_o] = 0$$

$$-2E_{y \sim p(y|x)}[(y - w_*^T x)x | x_o] = 0$$

By definition, we stated earlier that w_* is the optimal value of w that sends the predictive error to 0, so we see that the mean of the predictive errors is equal to 0 for $\hat{w} = w_*$. Thus, if we take any linear function of the input x , by linearity of expectation, $a^T x$ would still be uncorrelated to the predictive error. Thus,

$$E_{y \sim p(y|x)}[(y - w_*^T x)a^T x] = 0$$

Problem 2 [5 points] Explain (succinctly but precisely) how the original statement, regarding zero correlation between the prediction errors made by \hat{w} estimated by least squares with any linear function of the training $\{x_i\}$, follows from (3).

It follows from the result of 3 that the predictive error has a covariance of 0 with any linear function of x , since the definition of $Cov(U, V) = E[UV] - E[U]E[V]$. We know that if we model the predictive error as a random variable U , and linear functions of x as V , then we have shown already that $E[UV] = 0$.

Next, we consider the effect of scaling the input in the regression problem. This becomes relevant, for instance, in polynomial regression with high order models, to prevent very large or very small values and the resulting potential for numerical instability (think about the case of $x = 0.01$ or $x = 1000$ for 10-order model...). It seems that a good solution could be to multiply each column of the design matrix X by a suitable number so that the range of that column (i.e. of the corresponding dimension in the regression input space) be fixed, say, to $[-1, 1]$.

Suppose that the data are scaled by multiplying the j -th dimension of the input by a non-zero number c_j . We will denote a single normalized data point by $\tilde{x} \equiv [1, c_1x_1, \dots, c_dx_d]^\top$ and the resulting design matrix by \tilde{X} .

Problem 3 [20 points] Let $\hat{\mathbf{w}}$ be the least squares estimate of the regression parameters from the unscaled \mathbf{X} , and let $\tilde{\mathbf{w}}$ be the solution obtained from the scaled $\tilde{\mathbf{X}}$. Show that the scaling does not change optimality, in the sense that $\hat{\mathbf{w}} \cdot \mathbf{x} = \tilde{\mathbf{w}} \cdot \tilde{\mathbf{x}}$

Hint: You may find it helpful to express the scaling as a linear operator, yielding a matrix-product expression, and to equip yourself with a matrix reference, such as the Matrix Cookbook (look it up if you are not already familiar with it!)

Let's begin by defining the scaling operator as a square diagonal matrix such that the normalized data point $\tilde{x} = Cx$ and the entire normalized design matrix is equal to $\tilde{X} = XC$. Now, let us define the loss function for the new normalized least squares estimate $L(\tilde{w}, \tilde{X}, y) = \frac{1}{N}(y - \tilde{X}\tilde{w})^T(y - \tilde{X}\tilde{w}) = \frac{1}{N}(y^T - \tilde{w}^T \tilde{X}^T)(y - \tilde{X}\tilde{w})$. Now, as usual, differentiate with respect to w to find the optimal \tilde{w} that minimizes the loss function.

$$\frac{\partial L(\tilde{w}, \tilde{X}, y)}{\partial w} = \frac{1}{N} \frac{\partial}{\partial w} [y^T y - y^T \tilde{X} \tilde{w} - \tilde{w}^T \tilde{X}^T y + \tilde{w}^T \tilde{X}^T \tilde{X} \tilde{w}]$$

$$0 = \frac{1}{N} [0 - 2\tilde{X}^T y + 2\tilde{X}^T \tilde{X} \tilde{w}]$$

$$\tilde{X}^T \tilde{X} \tilde{w} = \tilde{X}^T y$$

$$\tilde{w} = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T y = [(XC)^T (XC)]^{-1} (XC)^T y$$

Since C is diagonal, it is also symmetric so $C^T = C$. Also, $X^T X$ is symmetric

$$\tilde{w} = (CX^T XC)^{-1} CX^T y = (C^2(X^T X))^{-1} CX^T y$$

$$\tilde{w} = C^{-1}(X^T X)^{-1} X^T y = C^{-1} \hat{w}$$

$$\tilde{w} \cdot \tilde{x} = \hat{w} \cdot x = C^{-1} \hat{w} \cdot Cx$$

$$C^{-1} \hat{w} \cdot Cx = \hat{w} \cdot x$$

$$w^T C^{-1} Cx = \hat{w} \cdot x$$

$$w^T X = \hat{w} \cdot x$$

When the dimensionality d of the data (aka the number of features) is much larger than the number of training instances n (i.e. $d \gg n$), the matrix $\mathbf{X}^\top \mathbf{X}$ is not full rank and thus can not be inverted. Therefore, instead of minimizing the least squares loss, we minimize the following loss function:

$$\hat{R}_\lambda(\mathbf{w}) = \sum_{i=1}^n (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 + \lambda \sum_{j=1}^d w_j^2 = (\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2. \quad (4)$$

Problem 4 [20 points] Find $\hat{\mathbf{w}}$ that minimizes the empirical risk in equation (4).

$$\widehat{R}_\lambda(w) = (y - Xw)^T(y - Xw) + \lambda\|w\|_2^2$$

$$\widehat{R}_\lambda(w) = (y^T - w^T X^T)(y - Xw) + \lambda\|w\|_2^2$$

$$\widehat{R}_\lambda(w) = y^T y - 2(X^T w^T y) + w^T X^T X w + \lambda\|w\|_2^2$$

$$\widehat{R}_\lambda(w) = y^T y - 2(X^T w^T y) + w^T X^T X w + \lambda\|w\|_2^2$$

$$\frac{\partial \widehat{R}_\lambda(w)}{\partial w} = 0 - 2(X^T y) + 2X^T X \widehat{w} + 2\lambda \widehat{w} = 0$$

$$0 = -(X^T y) + X^T X \widehat{w} + \lambda \widehat{w}$$

$$(X^T y) = X^T X \widehat{w} + \lambda \widehat{w}$$

$$(X^T y) = \widehat{w}(X^T X + \lambda)$$

$$(X^T y)(\lambda + X^T X)^{-1} = \widehat{w}$$

II. Asymmetric loss

In this section, we will consider a regression model associated with a different loss function.

In class we have discussed the squared loss

$$\ell(\hat{y}, y) = (\hat{y} - y)^2, \quad (5)$$

which is symmetric – the sign of the prediction error in by with respect to ground truth y doesn't matter. In some applications, this may not be a reasonable loss. For instance, consider the problem of estimating pressure necessary to put in the tire of a truck as a function of some parameters of the road, environment, vehicle etc. Putting pressure which is too low could reduce performance; if it's extremely low it can of course be dangerous. However, putting too much pressure may quickly result in catastrophic damage to the vehicle.

In this situation, we would like to penalize for positive errors differently (more severely) than for the negative errors of the magnitude. We can express this in the following, asymmetric squared loss:

$$\ell_{\alpha}(\hat{y}, y) = \begin{cases} (\hat{y} - y)^2 & \text{if } \hat{y} \leq y, \\ \alpha(\hat{y} - y)^2 & \text{if } \hat{y} \geq y. \end{cases} \quad (6)$$

The coefficient α determines how much more we worry about over-predicting y than about under-predicting. We will assume that $\alpha \geq 1$. See Figure 1 for an illustration.

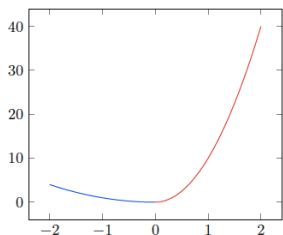


Figure 1: Plot of asymmetric loss ℓ_{10} ; blue portion corresponds to negative errors (under-prediction), red to positive errors.

For the following problems, we have provided you with three data sets: training, validation and test, drawn from the same (unknown to you) $p(x, y)$. The data are in a HDF5 file `data.h5`. You will use training set to fit different models; validation to select between models, and test set to report the accuracy with the model that you end up choosing. That means that the test set will be only used once, to compute and report test loss with the model of choice.

For your convenience, we provide most of the code you will need in the experiments below, with a few missing pieces that you will need to fill in. If you are new to Python, don't worry too much about efficiency and elegance of your code, although those are of course good goals eventually; for now just make sure it works correctly and is readable.

The code “skeleton”, with quite a bit of flesh on it, is provided as a Jupyter notebook `PS1-asymloss-notebook.ipynb`, along with an auxiliary file `utils.py` containing a couple of functions you will need; take time to read and understand the provided code, especially if you are new to Python.

Please submit the notebook with all the code filled in and with the output of your experiments, after the name change as stated in the preamble of this problem set. Please write your comments,

explanations and interpretation of the results as requested in the problems below, in the PDF solution file.

Problem 5 [15 points]

Fit linear, quadratic and cubic models to the training data under the standard (symmetric) squared loss (5). Visualize the fit under the three models. To make it easier to parse your figures, let's agree on the color scheme: the plots will be blue (linear), green (quadratic) and red (cubic) solid lines.

Evaluate, for each of the three models, the (symmetric) squared loss and the log-likelihood of the training data under the Gaussian noise assumption.

Note: you will need to derive ML estimate for the Gaussian noise variance σ^2 , in addition to the ML estimate for the model parameters w ; write the missing code snippet for this, and explain in the solution file how you derived it.

Finally, compute and report (symmetric) squared loss for each of the three models on the validation set. Comment on the relative performance of the models, and on the gap, if any, between validation and training loss values. Select one of the three models based on this evaluation, and explain your selection. Let's call the model you end up choosing here model A.

degree 1:
train loss 0.155770
val loss 0.108249
sigma²: 0.158948
log-likelihood -0.489351

degree 2:
train loss 0.032250
val loss 0.024095
sigma²: 0.032908
log-likelihood 0.298083

degree 3:
train loss 0.031714
val loss 0.026440
sigma²: 0.032362
log-likelihood 0.306453

Empirically, it was found that a model of degree 2 has the lowest validation error but slightly higher training loss and variance, and lower log-likelihood (a higher log-likelihood is better since maximizing log-likelihood minimizes log-loss). These numbers make sense since without cross validation, the model that is chosen will always be more complex since it will always have the lowest training loss without fail. However, the model of degree 2 has the lowest validation loss, meaning it performs the best in terms of generalizing out of the three, so we choose model A to be of degree 2. For the MLE for Gaussian noise variance, we know that the noise is a 0 mean Gaussian random variable, so the derivation process is simply, which ultimately just leads us to the fact that the estimator for Gaussian noise is the sample variance.

Problem 6 [25 points]

Now we want to repeat the experiment above but under the asymmetric loss function (6). Since there is no closed form solution, we will need to rely on gradient descent. Derive the expression for the gradient of the loss (6) with respect to the parameter vector w , for a general form of linear regression (i.e., with inputs represented by a feature vector $\phi(x)$, which could capture linear, quadratic or cubic model).

Complete the provided code for gradient descent using the derivation for the gradient, and run it to fit linear, quadratic and cubic models under asymmetric squared loss. Visualize the resulting three models (use same color scheme as before, but dashed lines) and report the asymmetric squared loss values and the log-likelihood of the model *under the standard Gaussian noise model* on training data, and the asymmetric squared loss on validation data. Compare these numbers to those you obtained with symmetric loss, and comment on the relative performance of the three models, as well as the performance of the equal degree models across the two loss functions.

Next, using the outcome of your experiment, select a model among the three (linear / quadratic / cubic) and explain your selection. Let's call the resulting model B.

Evaluate model A (selected with symmetric loss) and model B on test set, and interpret your findings; which model do you conclude to be a better choice based on this entire experiment? Choice of evaluation metrics is up to you, and you should motivate it succinctly but clearly.

Hint: You may want to start by reviewing the derivation of the gradient of squared loss, and extend it to asymmetric loss.

degree 1:
train loss 0.456437
val loss 0.374340
sigma²: 0.159223
log-likelihood -1.005706

degree 2:
train loss 0.089027
val loss 0.076941
sigma²: 0.036509
log-likelihood -0.122804

degree 3:
train loss 0.088523
val loss 0.078923
sigma²: 0.036320
log-likelihood -0.114172

Here, based on a set of the same training data, we found that a degree 2 polynomial using an asymmetric loss model yields the lowest validation error, so we should use this as our model for the entire experiment.