

Homework #3

CSE 546: Machine Learning

Michael Ross

1 Bayesian Inference

1. [5 points] Let $\{(x_i, y_i)\}_{i=1}^n$ be sampled iid from a joint distribution P_{XY} over $\mathbb{R}^d \times \mathbb{R}$ such that for some $w \in \mathbb{R}^d$ we have $y_i \sim \mathcal{N}(x_i^T w, \sigma^2)$. That is, $p(Y = y|x, w) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{(y-x^T w)^2}{2\sigma^2})$. Express your answers in terms of $\mathbf{X} = [x_1, \dots, x_n]^T$ and $\mathbf{y} = [y_1, \dots, y_n]^T$.

- a. If $(\mathbf{X}^T \mathbf{X})^{-1}$ exists, what is the MLE of w ?

Answer:

Since y is drawn from a Gaussian, the MLE of w is the least square estimate of w :

$$\hat{w}_{MLE} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- b. Assume that the true w is drawn from a Gaussian prior distribution $p(w) = \frac{1}{(2\pi\tau^2)^{d/2}} \exp(-\frac{\|w\|_2^2}{2\tau^2})$. What is the MAP estimate of w ?

Answer:

$$\begin{aligned} \hat{w} &= \arg \max_w p(w|\mathbf{X}, \mathbf{y}) \\ &= \arg \max_w [p(\mathbf{y}|w, \mathbf{X})p(w)] \\ &= \arg \max_w \left[\frac{1}{(2\pi\sigma^2)^{d/2}} \exp(-\frac{\|\mathbf{y}-\mathbf{X}w\|_2^2}{2\sigma^2}) \frac{1}{(2\pi\tau^2)^{d/2}} \exp(-\frac{\|w\|_2^2}{2\tau^2}) \right] \\ &= \arg \max_w \left[\exp(-\frac{\|\mathbf{y}-\mathbf{X}w\|_2^2}{2\sigma^2} - \frac{\|w\|_2^2}{2\tau^2}) \right] \end{aligned}$$

Taking the derivative and setting equal to zero:

$$0 = \exp(-\frac{\|\mathbf{y}-\mathbf{X}w\|_2^2}{2\sigma^2} - \frac{\|w\|_2^2}{2\tau^2}) \left(\frac{\mathbf{X}^T(\mathbf{y}-\mathbf{X}w)}{\sigma^2} - \frac{w}{\tau^2} \right)$$

$$\frac{\mathbf{X}^T(\mathbf{y}-\mathbf{X}w)}{\sigma^2} = \frac{w}{\tau^2}$$

$$\mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{X} w = \frac{\sigma^2}{\tau^2} w$$

$$\hat{w}_{MAP} = (\frac{\sigma^2}{\tau^2} \mathbf{I} + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- c. Assuming the setting of part b, what is the posterior distribution $p(w|\mathbf{X}, \mathbf{y})$ of w ? Give your answer in terms of $\mathcal{N}(\mu, \Sigma)$ for some μ, Σ . What is $\mathbb{E}[w|\mathbf{X}, \mathbf{y}]$?

Answer:

$$\begin{aligned} p(w|\mathbf{X}, \mathbf{y}) &\propto p(\mathbf{y}|w, \mathbf{X})p(w) \\ &\propto \exp(-\frac{\|\mathbf{y}-\mathbf{X}w\|_2^2}{2\sigma^2} - \frac{\|w\|_2^2}{2\tau^2}) \\ &\propto \exp(-\frac{(\mathbf{y}-\mathbf{X}w)^T(\mathbf{y}-\mathbf{X}w)}{2\sigma^2} - \frac{w^T w}{2\tau^2}) \\ &\propto \exp(-\frac{(\mathbf{y}^T \mathbf{y} - w^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} w + w^T \mathbf{X}^T \mathbf{X} w)}{2\sigma^2} - \frac{w^T w}{2\tau^2}) \end{aligned}$$

Since $w^T \mathbf{X}^T \mathbf{y}$ is a constant, $w^T \mathbf{X}^T \mathbf{y} = \mathbf{y}^T \mathbf{X} w$

Dropping term constant with respect to w

$$\begin{aligned} &\propto \exp(-\frac{1}{2\sigma^2} (w^T \mathbf{X}^T \mathbf{X} w - 2w^T \mathbf{X}^T \mathbf{y} + \frac{\sigma^2}{\tau^2} w^T w)) \\ &\propto \exp(-\frac{1}{2\sigma^2} (w^T (\mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{I}) w - 2w^T \mathbf{X}^T \mathbf{y})) \\ &\propto \exp(-\frac{(\mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{I})}{2\sigma^2} (w^T w - 2w^T (\mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y})) \end{aligned}$$

Adding constant to complete the square

$$\begin{aligned} &\propto \exp\left(-\frac{(\mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{I})}{2\sigma^2} \|w - (\mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}\|_2^2\right) \\ &= \mathcal{N}\left((\mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}, \sigma^2 (\mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{I})^{-1}\right) \end{aligned}$$

$$\mathbb{E}[w|\mathbf{X}, \mathbf{y}] = (\mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} = \hat{w}_{MAP}$$

- d. Fix a $z \in \mathbb{R}^d$. If $f_z = z^T w$ is the predicted function value at z . Show that

$$f_z|\mathbf{X}, \mathbf{y} \sim \mathcal{N}(z^T (\frac{\sigma^2}{\tau^2} I + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}, \sigma^2 z^T (\frac{\sigma^2}{\tau^2} I + \mathbf{X}^T \mathbf{X})^{-1} z)$$

Answer:

For a generic $x \sim \mathcal{N}(\mu, \Sigma)$ and $y = a^T x + b$
 $y \sim \mathcal{N}(a^T \mu + b, a^T \Sigma a)$

$$\begin{aligned} \text{Since } w &\sim \mathcal{N}((\mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}, \sigma^2 (\mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{I})^{-1}) \\ f_z &\sim \mathcal{N}(z^T (\mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}, z^T \sigma^2 (\mathbf{X}^T \mathbf{X} + \frac{\sigma^2}{\tau^2} \mathbf{I})^{-1} z) \end{aligned}$$

- e. The matrix inversion identity says that for matrices A, U, C, V of the appropriate sizes and when A^{-1} exists, we have

$$(A + UCV)^{-1} = A^{-1} - A^{-1}U(C^{-1} + VA^{-1}U)^{-1}VA^{-1}.$$

Use this identity to show that

$$f_z|\mathbf{X}, \mathbf{y} \sim \mathcal{N}(\mathbf{k}_z^T (\frac{\sigma^2}{\tau^2} I + \mathbf{K})^{-1} \mathbf{y}, \tau^2 \mathbf{k}_{zz} - \tau^2 \mathbf{k}_z^T (\frac{\sigma^2}{\tau^2} I + \mathbf{K})^{-1} \mathbf{k}_z)$$

where $\mathbf{K} = \mathbf{X}\mathbf{X}^T$, $\mathbf{k}_z = \mathbf{X}z$, and $\mathbf{k}_{zz} = z^T z$. How does the MAP estimate of f_z relate to the solution of Kernel ridge regression with a linear kernel evaluated at z ?

You have just derived what is known as Gaussian process regression. For more information, consult Rasmussen and Williams' *Gaussian Processes for Machine Learning* book: <http://www.gaussianprocess.org/>.

2. [1 points] Let $\{(x_i, y_i)\}_{i=1}^n$ be sampled iid from a joint distribution P_{XY} over $\mathbb{R}^d \times \mathbb{R}$ such that for some $w \in \mathbb{R}^d$ we have $y_i \sim \mathcal{N}(x_i^T w, \sigma^2)$. That is, $p(Y = y|x, w) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{(y-x^T w)^2}{2\sigma^2})$. Express your answers in terms of $\mathbf{X} = [x_1, \dots, x_n]^T$ and $\mathbf{y} = [y_1, \dots, y_n]^T$.

- Assume that the true w is drawn from a Laplace prior distribution $p(w) = \frac{1}{(2a)^d} \exp(-\frac{\|w\|_1}{a})$. What the MAP estimate of w ?
- The Laplace prior is not conjugate to the the normal likelihood. Is $\mathbb{E}[w|\mathbf{X}, \mathbf{y}]$ necessarily the same as the MAP estimate? If not, provide an example.

2 Kernel Regression

3. [6 points] First let's generate some data. Let $n = 30$ and $f(x) = 4 \sin(\pi x) \cos(6\pi x^2)$. For $i = 1, \dots, n$ let each x_i be drawn uniformly at random on $[0, 1]$ and $y_i = f(x_i) + \epsilon_i$ where $\epsilon_i \sim \mathcal{N}(0, 1)$. Using kernel ridge regression, build a predictor

$$\hat{\alpha} = \min_{\alpha} \|K\alpha - y\|^2 + \lambda \alpha^T K \alpha, \quad \hat{f}(x) = \sum_{i=1}^n \hat{\alpha}_i k(x_i, x)$$

where $K_{i,j} = k(x_i, x_j)$ is a kernel evaluation and λ is the regularization constant.

- a. Using leave-one-out cross validation, find a good λ and hyperparameter settings for the following kernels:

- $k_{poly}(x, z) = (1 + x^T z)^d$ where $d \in \mathbb{N}$ is a hyperparameter,

- $k_{rbf}(x, z) = \exp(-\gamma\|x - z\|^2)$ where $\gamma > 0$ is a hyperparameter¹.

Report the values of d , γ , and the λ values for both kernels.

- For a single plot per kernel, plot the original data $\{(x_i, y_i)\}_{i=1}^n$, the true $f(x)$, the $\hat{f}(x)$ found through leave-one-out CV.
- Using the fixed hyperparameters you found in part a, we wish to build Bootstrap percentile confidence intervals for $\hat{f}_{poly}(x)$ and $\hat{f}_{rbf}(x)$ for all $x \in [0, 1]$. Use the non-parametric bootstrap with $B = 300$ datasets (i.e. randomly draw with replacement n samples from $\{(x_i, y_i)\}_{i=1}^n$ and train an \hat{f} , repeat this B times) and find 5% and 95% percentiles (see Hastie, Tibshirani, Friedman Ch. 8.2 for a review). Plot the percentile curves on the plots from part b.
- Repeat all parts of this problem with $n = 300$ (you may just use 10-fold CV instead of leave-one-out)
- Suppose m additional samples are drawn i.i.d. the same way the first n samples were drawn. Propose a statistical significance test to decide which learned function (which kernel) is the better fit (hint: if $\epsilon_i \sim \mathcal{N}(0, 1)$, how is $\sum_i \epsilon_i^2$ distributed?).

3 k -means clustering

- [5 points] Given a dataset $x_1, \dots, x_n \in \mathbb{R}^d$ and an integer $1 \leq k \leq n$, recall the following k -means objective function

$$\min_{\pi_1, \dots, \pi_k} \sum_{i=1}^k \sum_{j \in \pi_i} \|x_j - \mu_i\|_2^2, \quad \mu_i = \frac{1}{|\pi_i|} \sum_{j \in \pi_i} x_j. \quad (1)$$

Above, $\{\pi_i\}_{i=1}^k$ is a partition of $\{1, 2, \dots, n\}$. The objective (1) is NP-hard² to find a global minimizer of. Nevertheless the commonly used heuristic which we discussed in lecture, known as Lloyd's algorithm, typically works well in practice. Implement Lloyd's algorithm for solving the k -means objective (1). Do not use any off the shelf implementations, such as those found in `scikit-learn`.

- Run the algorithm on MNIST with $k = 5, 10, 20$, plotting the objective function (1) as a function of iteration. Visualize (and include in your report) the cluster centers as a 28×28 image.
- Implement the `kmeans++` initialization scheme³ for your k -means implementation and repeat part a. Note that this initialization scheme is widely used in practice, and as a rule should be used. Plot the objective function as a function of iteration. Are the identified centers visually better than part a?

4 Joke Recommender System

- [8 points] You will build a personalized joke recommender system. There are $m = 100$ jokes and $n = 24,983$ users⁴. As historical data, every user read a subset of jokes and rated them. The goal is to recommend more jokes, such that the recommended jokes match the individual user's sense of humor. The historical rating is represented by a matrix $R \in \mathbb{R}^{n \times m}$. The entry $R_{i,j}$ represents the user i 's rating on joke j . The rating is a real number in $[-10, 10]$: a higher value represents that the user is more satisfied with the joke. The directory `/jokes` contains the text of all 100 jokes. Read them before you start! In addition, you are provided with two files:

- `train.txt` contains the joke-user-score data representing the training set. Each line takes the form "`i, j, s`", where `i` is the user index, `j` is the joke index, and `s` is the user's score in $[-10, 10]$ describing how much they liked the joke (higher is better).

¹Given a dataset $x_1, \dots, x_n \in \mathbb{R}^d$, a heuristic for choosing γ is the inverse of the median of all $\binom{n}{2}$ squared distances $\|x_i - x_j\|_2^2$.

²To be more precise, it is both NP-hard in d when $k = 2$ and k when $d = 2$. See the references on the wikipedia page for k -means for more details.

³See <http://ilpubs.stanford.edu:8090/778/1/2006-13.pdf>.

⁴Data from <http://eigentaste.berkeley.edu/dataset/>

- `test.txt` has the same format, with the same users rating movies held out from the training set.

Latent factor model is the state-of-the-art method for personalized recommendation. It learns a vector representation $u_i \in \mathbb{R}^d$ for each user and a vector representation $v_j \in \mathbb{R}^d$ for each joke, such that the inner product $\langle u_i, v_j \rangle$ approximates the rating $R_{i,j}$. You will build a simple latent factor model. We will evaluate our learnt vector representations by two metrics

- Mean squared error: $\frac{1}{|S|} \sum_{(i,j) \in S} (\langle u_i, v_j \rangle - R_{i,j})^2$ where S (and the corresponding $R_{i,j}$ values) are from the test set
- Mean absolute error: $\frac{1}{n} \sum_{i=1}^n \frac{1}{|\mathcal{N}_i|} \sum_{j \in \mathcal{N}_i} |\langle u_i, v_j \rangle - R_{i,j}|$ where \mathcal{N}_i are the jokes rated by user i in the test set

You will implement multiple estimators and use the inner product $\langle u_i, v_j \rangle$ to predict if user i likes joke j in the test data. You will choose hyperparameters like d or the amount of regularization by creating a validation set from the training set.

- The first estimator pools all the users together and just predicts what the average user in the training set rated the joke. This is equivalent to $d = 1$ with u as the all ones vector and v minimizing least squares.
- Now replace all missing values in $R_{i,j}$ no in the training set by zero. Then use singular value decomposition (SVD) to learn a lower dimensional vector representation for users and jokes. Recall this means to project the data vectors to lower dimensional subspaces of their corresponding spaces, spanned by singular vectors. Refer to the lecture materials on SVD, PCA and dimensionality reduction. You should use an efficient solver, I recommend `scipy.sparse.linalg.svds`. Try $d = 1, 2, 5, 10, 20, 50$ and plot the error metrics on the train and test as a function of d .
- For sparse data, replacing all missing values by zero is not a completely satisfying solution. A missing value means that the user has not read the joke, but doesn't mean that the rating should be zero. A more reasonable choice is to minimize the MSE only on rated jokes. Let's define a loss function:

$$L(\{u_i\}, \{v_j\}) := \sum_{(i,j) \in T} (\langle u_i, v_j \rangle - R_{i,j})^2 + \lambda \sum_{i=1}^n \|u_i\|_2^2 + \lambda \sum_{j=1}^m \|v_j\|_2^2,$$

where T and $R_{i,j}$ here are from the training set and $\lambda > 0$ is the regularization coefficient. Implement an algorithm to learn vector representations by minimizing the loss function $L(\{u_i\}, \{v_j\})$. Try $d = 1, 2, 5, 10, 20, 50$ and plot the error metrics on the train and test as a function of d . Note that you may need to tune the hyper-parameter λ to optimize the performance.

Hint: you may want to employ an alternating minimization scheme. First, randomly initialize $\{u_i\}$ and $\{v_j\}$. Then minimize the loss function with respect to $\{u_i\}$ by treating $\{v_j\}$ as constant vectors, and minimize the loss function with respect to $\{v_j\}$ by treating $\{u_i\}$ as constant vectors. Iterate these two steps until both $\{u_i\}$ and $\{v_j\}$ converge. Note that when one of $\{u_i\}$ or $\{v_j\}$ is given, minimizing the loss function with respect to the other part has closed-form solutions. You should never be allocating an $m \times n$ matrix for this problem.