

Homework 1

FINAL VERSION

Before this final version on January 20th, there was a version uploaded on January 19th. To avoid confusion, here we list the changes made from the previous version.

- Score distribution has been changed. The course staff decided that the problem 2 is much harder than problem 1, and downgraded the point of problem 1 to 30 points.
- Previous section 1.1 Expectation has been removed and now section 1.1 is Independence.
- Some more explanations have been added to problem 2.
- The deadline has been extended for a day from January 26th to January 27th.
- Lastly, the following instructions on next section were clarified.

INSTRUCTIONS

- The homework is due at 9:00am on January 27, 2015. Anything that is received after that time will be considered to be late.
- Answers to every theory questions need to be submitted electronically on Autolab. Only PDF is acceptable (e.g. LaTeX or handwritten and scanned).
- Make sure you prepare the answers to each question separately. This helps us dispatch the problems to different graders.
- Collaboration on solving the homework is allowed. Discussions are encouraged but you should think about the problems on your own.
- When you do collaborate, you should list your collaborators! Also cite your resources, in case you got some inspiration from other resources (books, websites, papers).
- If you do collaborate with someone or use a book or website, you are expected to write up your solution independently. That is, close the book and all of your notes before starting to write up your solution.

1 Probability Review and Bayesian Spam Filter [Zhou; 30pts]

Probability is, in many ways, the most fundamental mathematical technique for machine learning. This problem will review several basic notions from probability and make sure that you remember how to do some elementary proofs.

Recall that for a discrete random variable X whose values are integers, we frequently use the notation $P(X = x)$. If a random variable Y is continuous, we typically use a “density function” $f(Y = y)$. The conditions for $P(X = x)$ to be a valid probability distribution are that $\sum_{-\infty}^{\infty} P(X = x) = 1$ and $P(X = x) \geq 0$. Similarly for $f(Y = y)$ to be a valid continuous distribution, $\int_{-\infty}^{\infty} f(Y = y) dy = 1$ and $f(Y = y) \geq 0$.

Sometimes the underlying probability space has more than one variable (for example, the height and weight of a person). In this case, we may use notation like $f(X = x; Y = y)$ to denote the probability density function in several dimensions.

1.1 Independence [5 pts]

Intuitively, two random variables X and Y are “independent” if knowledge of the value of one tells you nothing at all about the value of the other. Precisely, if X and Y are discrete, independence means that $P(X = x; Y = y) = P(X = x)P(Y = y)$, and if they are continuous, $f(X = x; Y = y) = f(X = x)f(Y = y)$.

Homework 1

Show the following, for independent random variable X and Y :

$$E[XY] = E[X]E[Y] \text{ for discrete and continuous cases respectively.}$$

Answer:

Discrete:

$$\begin{aligned} E[XY] &= \sum_{X=x} \sum_{Y=y} xy P(X=x, Y=y) \\ &= \sum_{X=x} x \sum_{Y=y} y P(X=x) P(Y=y) \\ &= \sum_{X=x} x P(X=x) (\sum_{Y=y} y P(Y=y)) \\ &= \sum_{X=x} x P(X=x) E[Y] \\ &= E[Y] \sum_{X=x} x P(X=x) \\ &= E[Y] E[X] \end{aligned}$$

Continuous:

$$\begin{aligned} E[XY] &= \int_x \int_y xy f_{XY}(x, y) dy dx \\ &= \int_x \int_y xy f_X(x) f_Y(y) dy dx \\ &= \int_x f_X(x) \left(\int_y y f_Y(y) dy \right) dx \\ &= \int_x x f_X(x) (E[Y]) dx \\ &= E[Y] \int_x x f_X(x) dx \\ &= E[Y] E[X] \end{aligned}$$

1.2 Spam filtering equation [5 pts]

Naive Bayes classifiers work by correlating the use of tokens (typically words, or sometimes other things), with spam and non-spam e-mails and then using Bayesian inference to calculate a probability that an email is or is not spam. Naive Bayes spam filtering is a baseline technique for dealing with spam. Let's suppose the suspected message contains the word "replica". Most people who are used to receiving e-mail know that this message is likely to be spam, more precisely a proposal to sell counterfeit copies of well-known brands of watches. The spam detection software, however, does not "know" such facts; all it can do is compute probabilities.

Let $Pr(S|W)$ be the probability that a message is a spam given the word "replica" appears in it. Express $Pr(S|W)$ in terms of the following components we provide:

1. $Pr(S)$ is the overall probability that any given message is spam;
2. $Pr(W|S)$ is the probability that the word "replica" appears in spam messages;
3. $Pr(H)$ is the overall probability that any given message is not spam (is "ham");
4. $Pr(W|H)$ is the probability that the word "replica" appears in ham messages.

Homework 1

Answer:

$$Pr(S|W) = \frac{Pr(W|S)Pr(S)}{Pr(W|S)Pr(S) + Pr(W|H)Pr(H)} \quad (1)$$

1.3 I.I.D. assumption in spam filters [10 pts]

Give at least 4 cases how "Independent and identically distributed random variables (i.i.d.)" assumption can be violated for spam filtering. For each case, please explain why i.i.d. is violated and how spammers might exploit the situation.

Answer:

- Rare words and symbols are not equally distributed, this is not reflected in the equation
- The relationship of the adjacent words are not modeled
- Word orders are not considered in the model
- Syntactic structures are not considered in the model

Note: answers are not limited for the above four.

1.4 Poison the Bayesian spam filter [10 pts]

Depending on the implementation, Bayesian spam filtering may be susceptible to Bayesian poisoning, a technique used by spammers in an attempt to degrade the effectiveness of spam filters that rely on Bayesian filtering. Think of some ideas to make your spam pass the Bayesian spam filter based on the equation.

Answer:

Spammer tactics include insertion of random innocuous words that are not normally associated with spam, thereby decreasing the email's spam score, making it more likely to slip past a Bayesian spam filter. However with (for example) Paul Graham's scheme only the most significant probabilities are used, so that padding the text out with non-spam-related words does not affect the detection probability significantly. Words that normally appear in large quantities in spam may also be transformed by spammers. For example, Viagra would be replaced with Viaagra or V!agra in the spam message. The recipient of the message can still read the changed words, but each of these words is met more rarely by the Bayesian filter, which hinders its learning process. As a general rule, this spamming technique does not work very well, because the derived words end up recognized by the filter just like the normal ones. Another technique used to try to defeat Bayesian spam filters is to replace text with pictures, either directly included or linked. The whole text of the message, or some part of it, is replaced with a picture where the same text is "drawn". The spam filter is usually unable to analyze this picture, which would contain the sensitive words like Viagra.

2 Regression [Jay-Yoon; 70pts]

The objective of this problem is to gain knowledge on Linear Regression, Maximum Likelihood Estimation (MLE), Maximum-a-Posteriori Estimation (MAP) and the variants of Regression problems with introduction of regularization terms.

Homework 1

2.1 Linear Regression: MLE and Least Squares [25 pts]

Consider a linear model with some Gaussian noise:

$$Y_i = \langle X_i, w \rangle + b + \epsilon_i \quad \text{where} \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2), \quad i = 1, \dots, n. \quad (2)$$

Where $Y_i \in \mathbb{R}$ is a scalar, $X_i \in \mathbb{R}^d$ is a d -dimensional vector, $b \in \mathbb{R}$ is a constant, $w \in \mathbb{R}^d$ is d -dimensional weight on X_i , and ϵ_i is a i.i.d. Gaussian noise with variance σ^2 . Given the data $X_i, i = 1, \dots, n$, it is our goal to estimate w and b which specify the model.

We will show that solving the linear model (2) with MLE method is the same as solving the following Least Squares problem,

$$\operatorname{argmin}_{\beta} (\mathbf{y} - \mathbf{X}'\beta)^\top (\mathbf{y} - \mathbf{X}'\beta), \quad (3)$$

where $Y = (Y_1, \dots, Y_n)^\top$, $X'_i = (1, X_i)^\top$, $\mathbf{X}' = (X'_1, \dots, X'_n)^\top$ and $\beta = (b, w)^\top$.

1. From the model (2), derive the conditional distribution of $Y_i|X_i, w, b$. Again, X_i is a fixed data point.

Answer:

Note that $Y_i|X_i, w, b \sim \mathcal{N}(\langle X_i, w \rangle + b, \sigma^2)$, thus we have pdf of $Y_i|X_i, w, b$ in the following form:

$$f(Y_i = y_i|X_i, w, b) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y_i - \langle X_i, w \rangle - b)^2}{2\sigma^2}}.$$

2. Assuming i.i.d. between each $\epsilon_i, i = 1, \dots, n$, give an explicit expression for the loglikelihood, $\ell(Y|\beta)$ of the data.

Note: Notation Y and β is given on (3). $X_i, i = 1, \dots, n$ is fixed data points. Also given ϵ_i s are i.i.d., so $P(Y|\beta) = \prod_i P(Y_i|w, b)$.

Answer:

Changed the notation of loglikelihood to ℓ and notation of pdf to $f(\cdot)$ as some of the students were confused. And also note that we are just omitting X_i for convenience, as the problem explicitly tells that X_i are fixed points. Given $\mathbf{y} = (y_1, \dots, y_n)^\top$, since Y_i are independent as ϵ_i s are i.i.d. and X_i are given, the likelihood of $Y|\beta$ is as follows:

$$\begin{aligned} f(\mathbf{Y} = \mathbf{y}|\beta) &= \prod_{i=1}^n f(y_i|w, b) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y_i - \langle X_i, w \rangle - b)^2}{2\sigma^2}} \\ &= \left(\frac{1}{\sqrt{2\pi}\sigma} \right)^n e^{-\frac{\sum_{i=1}^n (y_i - \langle X_i, w \rangle - b)^2}{2\sigma^2}}. \end{aligned}$$

Now, taking the log, loglikelihood of $Y|\beta$ is as follows:

$$\ell(Y = \mathbf{y}|\beta) = -n \log \sqrt{2\pi}\sigma - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \langle X_i, w \rangle - b)^2 \quad (4)$$

Homework 1

3. Now show that solving for β that maximizes the loglikelihood, i.e. MLE, is the same as solving the Least Square problem of (3).

Answer:

To maximize the loglikelihood $\ell(Y = \mathbf{y}|\beta)$, we want to focus on the second term as the first term of the loglikelihood (4) is a constant. In short, to maximize the second term of (4), we want to minimize the $\sum_{i=1}^n (y_i - \langle X_i, w \rangle - b)^2$. Writing in the matrix-vector form:

$$\begin{aligned} \max_{\beta} \ell(Y = \mathbf{y}|\beta) &= \min_{\beta} \sum_{i=1}^n (y_i - \langle X_i, w \rangle - b)^2 \\ &= \min_{\beta} (\mathbf{y} - \mathbf{X}'\beta)^\top (\mathbf{y} - \mathbf{X}'\beta), \end{aligned}$$

where again $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$, $X'_i = (1, X_i)^\top$, $\mathbf{X}' = (X'_1, \dots, X'_n)^\top$ and $\beta = (b, w)^\top$.

4. Derive β that maximizes the loglikelihood.

Assume \mathbf{X}' has full rank on column space.

Hint: This link has some reference to matrix calculus: <http://www.cs.nyu.edu/~roweis/notes/matrixid.pdf>.

Answer:

Setting the objective function $J(\beta)$ as following:

$$\begin{aligned} J(\beta) &\equiv (\mathbf{y} - \mathbf{X}'\beta)^\top (\mathbf{y} - \mathbf{X}'\beta) \\ \nabla_{\beta} J(\beta) &= 2\mathbf{X}'^\top (\mathbf{X}'\beta - \mathbf{y}), \end{aligned}$$

the loglikelihood maximizer $\hat{\beta}$ can be found by solving the following optimality condition:

$$\begin{aligned} \nabla_{\beta} J(\hat{\beta}) &= 0 \\ \iff \mathbf{X}'^\top (\mathbf{X}'\hat{\beta} - \mathbf{y}) &= 0 \\ \iff \mathbf{X}'^\top \mathbf{X}'\hat{\beta} &= \mathbf{X}'^\top \mathbf{y} \\ \iff \hat{\beta} &= (\mathbf{X}'^\top \mathbf{X}')^{-1} \mathbf{X}'^\top \mathbf{y} \end{aligned}$$

Note that $(\mathbf{X}'^\top \mathbf{X}')^{-1}$ is possible because \mathbf{X}' has full rank on column space.

2.2 Nonlinear Regression and Regularization [15 pts]

Consider higher order term, $\phi(X_i) = (1, X_i, X_i^2, \dots, X_i^k)^\top$, then we can model Y in k th-order model as following:

$$Y_i = \langle \phi(X_i), \beta' \rangle + \epsilon_i \quad \text{where} \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2), \quad i = 1, \dots, n. \quad (5)$$

where all the definitions are from equation (2) except $\beta' \in \mathbb{R}^{k+1}$ and for simplicity let $d = 1$ for $X_i \in \mathbb{R}^d$.

1. As we did in section 2.1, show that optimal β' of (5) is $(\phi(\mathbf{X})^\top \phi(\mathbf{X}))^{-1} \phi(\mathbf{X})^\top \mathbf{Y}$ if we use MLE.

Homework 1

Note: $\phi(\mathbf{X}) = (\phi(X_1), \phi(X_2), \dots, \phi(X_n))^T$ and assume $\phi(\mathbf{X})$ has full rank on column space.

Hint: You are not expected to write the whole steps again. Focus on the change from the loglikelihood expressions of 2.1 and derive the optimization problem.

answer:

As the relation between Y_i and X_i have changed, the conditional pdf and likelihood function changes as well. The expressions can be just modified by just replacing X_i to $\phi(X_i)$ as following:

$$f(Y_i = y_i | X_i, w, b) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y_i - \langle \phi(X_i), \beta' \rangle)^2}{2\sigma^2}}$$

$$\ell(Y = \mathbf{y} | \beta') = -n \log \sqrt{2\pi}\sigma - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \langle \phi(X_i), \beta' \rangle)^2.$$

As we did before the maximum likelihood method can be expressed as

$$\begin{aligned} \max_{\beta'} \ell(Y = \mathbf{y} | \beta') &= \min_{\beta'} \sum_{i=1}^n (y_i - \langle \phi(X_i), \beta' \rangle)^2 \\ &= \min_{\beta'} (\mathbf{y} - \phi(\mathbf{X})\beta')^\top (\mathbf{y} - \phi(\mathbf{X})\beta'), \end{aligned}$$

and thus the maximum likelihood estimator $\hat{\beta} = (\phi(\mathbf{X})^\top \phi(\mathbf{X}))^{-1} \phi(\mathbf{X})^\top \mathbf{y}$ through similar steps as 2.1.

2. In case $\phi(\mathbf{X})^\top \phi(\mathbf{X})$ is not invertible, you can add diagonal term $\lambda \mathbf{I}_{k+1}$ to it so that $\phi(\mathbf{X})^\top \phi(\mathbf{X}) + \lambda \mathbf{I}_{k+1}$ becomes invertible with \mathbf{I}_{k+1} , an identity matrix of size n , and $\lambda > 0$. Show that $(\phi(\mathbf{X})^\top \phi(\mathbf{X}) + \lambda \mathbf{I}_{k+1})^{-1} \phi(\mathbf{X})^\top \mathbf{Y}$ is the solution of the optimization problem,

$$\operatorname{argmin}_{\beta'} \|\mathbf{Y} - \phi(\mathbf{X})\beta'\|_2^2 + \lambda \|\beta'\|_2^2. \quad (6)$$

answer:

The procedure is almost the same with 2.1.4. First, set the objective function $J'(\beta)$ as following:

$$\begin{aligned} J'(\beta) &\equiv (\mathbf{y} - \phi(\mathbf{X})\beta)^\top (\mathbf{y} - \phi(\mathbf{X})\beta) + \lambda \|\beta'\|_2^2 \\ \implies \nabla_{\beta'} J'(\beta') &= 2\phi(\mathbf{X}')^\top (\phi(\mathbf{X}')\beta' - \mathbf{y}) + 2\lambda\beta'. \end{aligned}$$

The loglikelihood maximizer $\hat{\beta}'$ can be found by solving the following optimality condition:

$$\begin{aligned} \nabla_{\beta'} J'(\hat{\beta}') &= 0 \\ \iff \phi(\mathbf{X}')^\top (\phi(\mathbf{X}')\hat{\beta}' - \mathbf{y}) + \lambda\hat{\beta}' &= 0 \\ \iff (\phi(\mathbf{X}')^\top \phi(\mathbf{X}') + \lambda \mathbf{I}_{k+1}) \hat{\beta}' &= \phi(\mathbf{X}')^\top \mathbf{y} \\ \iff \hat{\beta}' &= (\phi(\mathbf{X}')^\top \phi(\mathbf{X}') + \lambda \mathbf{I}_{k+1})^{-1} \phi(\mathbf{X}')^\top \mathbf{y} \end{aligned}$$

Homework 1

2.3 MAP [30 pts]

From (5), now consider a case where β' has a prior distribution $\beta' \sim \mathcal{N}(0, \eta^2 \mathbf{I}_{k+1})$.

1. Write the posterior distribution of $\beta'|Y_i$ given i th sample and $\beta'|\mathbf{Y}$ given whole data, respectively. Assume independence between β' and noise $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$, $i = 1, \dots, n$.

Hint1: Use Baye's rule and follow similar steps with 2.1. $\Pr(\beta'|Y_i) = \frac{\Pr(\beta')\Pr(Y_i|\beta')}{\Pr(Y_i)}$

Hint2: Sum of two independent Normal variable follows Normal distribution. i.e. $X \sim \mathcal{N}(a, b^2)$, $Y \sim \mathcal{N}(c, d^2)$, $X \perp Y \rightarrow X + Y \sim \mathcal{N}(a + c, b^2 + d^2)$.

answer:

From problem 2.2, we are aware that $Y_i|X_i, \beta' \sim \mathcal{N}(\langle \phi(X_i), \beta' \rangle, \sigma^2)$ and from this problem that $\beta' \sim \mathcal{N}(0, \eta^2 \mathbf{I}_{k+1})$. Using Baye's rule,

$$\begin{aligned} h(\beta'|Y_i) &\propto f(Y_i|\beta')g(\beta') \\ &\propto \exp\left(-\frac{1}{2\sigma^2}(Y_i - \langle \phi(X_i), \beta' \rangle)^2\right) \cdot \exp\left(-\frac{1}{2\eta^2}\beta'^T \beta'\right) \end{aligned}$$

where $g(\cdot)$ is the density function for β' and $h(\cdot)$ is the density function for $\beta'|Y_i$.

The normalizer for pdf h is defined as $Z = \int_{-\infty}^{\infty} f(Y_i|\beta')g(\beta')d\beta'$ and rewriting the pdf of $h(\beta'|Y_i)$,

$$h(\beta'|Y_i) = \frac{1}{Z} \exp\left(-\frac{1}{2\sigma^2}(Y_i - \langle \phi(X_i), \beta' \rangle)^2 - \frac{1}{2\eta^2}\beta'^T \beta'\right)$$

Likewise, pdf of $\beta'|\mathbf{Y}$ is

$$\begin{aligned} h(\beta'|\mathbf{Y}) &= \frac{1}{Z_n} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{Y} - \phi(\mathbf{X})\beta')^T (\mathbf{Y} - \phi(\mathbf{X})\beta') - \frac{1}{2\eta^2}\beta'^T \beta'\right) \\ &= \frac{1}{Z_n} \exp\left(-\frac{1}{2\sigma^2}\|\mathbf{Y} - \phi(\mathbf{X})\beta'\|_2^2 - \frac{1}{2\eta^2}\|\beta'\|_2^2\right) \end{aligned}$$

where Z_n is normalization factor defined as $Z_n = \int_{-\infty}^{\infty} f(\mathbf{Y}|\beta')g(\beta')d\beta'$.

2. MAP estimate β' is defined as β' that establishes the mode of posterior $P(\beta'|\mathbf{X}, \mathbf{Y})$. Show that solving for MAP estimate leads to the problem (6) if λ can be expressed in terms σ and η , i.e. $\lambda = g(\sigma, \eta)$. Find explicit expression for $g(\sigma, \eta)$.

answer: From the previous subproblem, it is clear that

$$\max_{\beta'} h(\beta'|Y_i) = \min_{\beta'} \frac{1}{2\sigma^2}\|\mathbf{Y} - \phi(\mathbf{X})\beta'\|_2^2 - \frac{1}{2\eta^2}\|\beta'\|_2^2. \quad (7)$$

If η becomes $\frac{\lambda}{\sigma^2}$, then the minimization problem becomes equivalent to (6),

$$\operatorname{argmin}_{\beta'} \|\mathbf{Y} - \phi(\mathbf{X})\beta'\|_2^2 + \lambda\|\beta'\|_2^2.$$

In other words, rearranging the terms, $\lambda = g(\sigma, \eta) = \frac{\sigma^2}{\eta^2}$ for (6) to become (7).

Homework 1

3. Describe one potential problem in the absence of regularization term in (6) and how the regularization term can alleviate the potential problem.

Note: Let's skip the case of non-invertible matrix as we already covered the case in section 2.2.

answer: The regularizing term penalizes large components in β' which leads to shrinking β' to have smaller norm. As a consequence, the penalty terms encourage the model to avoid overfitting and thus prevent adjusting to outlier data points which would otherwise influence β' drastically.