

Homework #6

TA in charge: Hanhsing Tu, Room 536

RELEASE DATE: 12/11/2008

DUE DATE: 12/18/2008, 4:00 pm IN CLASS

TA SESSION: 12/17/2008, noon to 2:00 pm IN R106

Unless granted by the instructor in advance, you must turn in a hard copy of your solutions (without the source code) for all problems. For problems marked with (), please follow the guidelines on the course website and upload your source code to designated places.*

Discussions on course materials and homework solutions are encouraged. But you should write the final solutions alone and understand them fully. Books, notes, and Internet resources can be consulted, but not copied from.

Since everyone needs to write the final solutions alone, there is absolutely no need to lend your homework solutions and/or source codes to your classmates at any time. In order to maximize the level of fairness in this class, lending and borrowing homework solutions are both regarded as dishonest behaviors and will be punished according to the honesty policy.

You should write your solutions in English with the common math notations introduced in class or in the problems. We do not accept solutions written in any other languages.

6.1 Bayesian Universe

(1) (15%) ASSUME that the universe generates an example (x, y) by the following procedure:

- (a) generate x from some probability density function $P(x)$
- (b) use some fixed (w, θ) to evaluate $\rho = \langle w, x \rangle - \theta$
- (c) generate $y \in \mathbb{R}$ from ρ by the probability density function $P(y|\rho) = \frac{1}{\sqrt{2\pi}} \exp(-(y - \rho)^2)$

If each (x_n, y_n) within $Z = \{(x_n, y_n)\}_{n=1}^N$ is generated i.i.d from the procedure above, what is the likelihood $P(Z|(w, \theta))$? Prove that linear regression (see Problem 2.3-(1)) equivalently gives the **maximum likelihood** estimate of (w, θ) .

(2) (20%) ASSUME that the universe generates an example (x, y) by the following procedure:

- (a) generate (w, θ) from $P(w, \theta) = \frac{1}{(\sqrt{2\pi})^{d+1} \cdot \sigma^{d+1}} \cdot \exp\left(-\frac{\|w\|^2 + \theta^2}{2\sigma^2}\right)$
- (b) generate x from some probability density function $P(x)$
- (c) use the “fixed” (w, θ) to evaluate $\rho = \langle w, x \rangle - \theta$
- (d) generate $y \in \mathbb{R}$ from ρ by the probability density function $P(y|\rho) = \frac{1}{\sqrt{2\pi}} \exp(-(y - \rho)^2)$

If each (x_n, y_n) within $Z = \{(x_n, y_n)\}_{n=1}^N$ is generated i.i.d from the procedure above, and assume that the constant $P(Z) = Q$, what is the posterior $P((w, \theta)|Z)$? Prove that regularized linear regression (see Problem 2.3.(3)) equivalently gives the **maximum a posteriori** estimate of (w, θ) . In particular, what is the relationship between λ (in Problem 2.3(3)) and σ (here)?

(3) (15%) ASSUME that the universe generates an example (x, y) by the following procedure:

- (a) generate x from some probability density function $P(x)$
- (b) use some fixed (w, θ) to evaluate $\rho = \langle w, x \rangle - \theta$
- (c) evaluate $Q_+ = \exp(\frac{\rho}{2})$ and $Q_- = \exp(-\frac{\rho}{2})$
- (d) generate $y \in \{+, -\}$ with the probability distribution $Q_y/(Q_+ + Q_-)$

If each (x_n, y_n) within $Z = \{(x_n, y_n)\}_{n=1}^N$ is generated i.i.d from the procedure above, what is the likelihood $P(Z|(w, \theta))$? Prove that logistic regression (see Problem 2.4) equivalently gives the **maximum likelihood** estimate of (w, θ) .

6.2 Power of Adaptive Boosting

The adaptive boosting (AdaBoost) algorithm, as shown in the class slides, is as follows:

- Input: $Z = \{(x_n, y_n)\}_{n=1}^N$.

- Set $u_n = \frac{1}{N}$ for all n .

- For $t = 1, 2, \dots, T$,

- Learn a simple rule h_t such that h_t solves

$$h_t = \min_h \sum_{n=1}^N u_n \cdot I[y_n \neq h(x_n)].$$

with the help of some base learner A_b .

- Compute the weighted error $\epsilon_t = \frac{1}{\sum_{m=1}^N u_m} \sum_{n=1}^N u_n \cdot I[y_n \neq h_t(x_n)]$ and the confidence

$$\alpha_t = \frac{1}{2} \ln \frac{1 - \epsilon_t}{\epsilon_t}$$

- Emphasize the training examples that do not agree with h_t :

$$u_n = u_n \cdot \exp(-\alpha_t y_n h_t(x_n)).$$

- Output: combined function $H(x) = \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(x) \right)$

In this problem, we will prove that AdaBoost can reach $\nu(H) = 0$ if T is large enough and every hypothesis h_t satisfies $\epsilon_t \leq \epsilon < \frac{1}{2}$.

- (1) (5%) Let $U^{(t-1)} = \sum_{n=1}^N u_n$ at the beginning of the t -th iteration. What is $U^{(0)}$?

- (2) (10%) According to the AdaBoost algorithm above, for $t \geq 1$, prove that

$$U^{(t)} = \frac{1}{N} \sum_{n=1}^N \exp \left(-y_n \sum_{\tau=1}^t \alpha_\tau h_\tau(x_n) \right).$$

- (3) (5%) By the result in (2), prove that $\nu(H) \leq U^{(T)}$.

- (4) (10%) According to the AdaBoost algorithm above, for $t \geq 1$, prove that

$$U^{(t)} = U^{(t-1)} \cdot 2\sqrt{\epsilon_t(1 - \epsilon_t)}.$$

- (5) (5%) Using $0 \leq \epsilon_t \leq \epsilon < \frac{1}{2}$, for $t \geq 1$, prove that

$$\sqrt{\epsilon_t(1 - \epsilon_t)} \leq \sqrt{\epsilon(1 - \epsilon)}.$$

- (6) (5%) Using $\epsilon < \frac{1}{2}$, prove that

$$\sqrt{\epsilon(1 - \epsilon)} \leq \frac{1}{2} \exp \left(-2 \left(\frac{1}{2} - \epsilon \right)^2 \right).$$

- (7) (5%) Using the results above, prove that

$$U^{(T)} \leq \exp \left(-2T \left(\frac{1}{2} - \epsilon \right)^2 \right).$$

- (8) (5%) Using the results above, argue that after $T = O(\log N)$ iterations, $\nu(H) = 0$.

6.3 Experiments with Bootstrap Aggregation (*)

- (1) (20%) Implement the decision stump learning algorithm A_{ds} . That is, let

$$h_{s,i,\theta}(x) = \text{sign}(s \cdot (x)_i - \theta),$$

where $s \in \{-1, +1\}$, $i \in \{1, 2, \dots, d\}$, and $\theta \in \mathbb{R}$. Given a weighted training set $Z = \{(x_n, y_n, u_n)\}_{n=1}^N$,

$$A_{ds}(Z) = \underset{h_{s,i,\theta}}{\text{argmin}} \sum_{n=1}^N u_n \cdot I[y_n \neq h_{s,i,\theta}(x_n)].$$

Run the algorithm on the following set for training (with $u_n = \frac{1}{N}$ for all N):

http://www.csie.ntu.edu.tw/~htlin/course/ml08fall/data/hw6_train.dat

and the following set for testing:

http://www.csie.ntu.edu.tw/~htlin/course/ml08fall/data/hw6_test.dat

Let g be the decision function returned from A_{ds} . Report $\nu(g)$ and $\hat{\pi}(g)$. Briefly state your findings.

- (2) (30%) Implement the bootstrap aggregation (bagging) algorithm with decision stumps (i.e., use A_{ds} as A_b below):

- Input: $Z = \{(x_n, y_n)\}_{n=1}^N$.
- for $t = 1, 2, \dots, T$,
 - generate $Z^{(t)}$ from Z by bootstrapping—uniformly sampling N examples from Z with replacement
 - let $h_t = A_b(Z^{(t)})$ and $\alpha_t = 1$.
- Output: combined function $H(x) = \text{sign}\left(\sum_{t=1}^T \alpha_t h_t(x)\right)$

Use a total of $T = 100$ iterations. Let $H_t(x) = \text{sign}\left(\sum_{\tau=1}^t \alpha_\tau h_\tau(x)\right)$. Plot $\nu(H_t)$ and $\hat{\pi}(H_t)$ as functions of t on the same figure. Briefly state your findings.

- (3) (Bonus 5%) Prove that you can implement an A_{ds} that runs in time $O(N \log N)$ instead of the brute-force implementation that takes $O(N^2)$.

6.4 Experiments with Adaptive Boosting (*)

Implement the AdaBoost algorithm (as in Problem 6.2 above) with decision stumps (i.e., use A_{ds} as A_b). Run the algorithm on the following set for training:

http://www.csie.ntu.edu.tw/~htlin/course/ml08fall/data/hw6_train.dat

and the following set for testing:

http://www.csie.ntu.edu.tw/~htlin/course/ml08fall/data/hw6_test.dat

- (1) (30%) Use a total of $T = 100$ iterations. Let $H_t(x) = \text{sign}\left(\sum_{\tau=1}^t \alpha_\tau h_\tau(x)\right)$. Plot $\nu(H_t)$, $\hat{\pi}(H_t)$, **AND** $U^{(t)}$ (see the definition above) as functions of t on the same figure. Briefly state your findings.
- (2) (20%) Compare your plots in Problem 6.3 and Problem 6.4. Briefly state your findings.