

Maximum Likelihood, Logistic Regression, and Stochastic Gradient Training

Charles Elkan
elkan@cs.ucsd.edu

January 10, 2012

1 Principle of maximum likelihood

Consider a family of probability distributions defined by a set of parameters θ . The distributions may be either probability mass functions (pmfs) or probability density functions (pdfs). Suppose that we have a random sample drawn from a fixed but unknown member of this family. The random sample is a training set of n examples x_1 to x_n . An example may also be called an observation, an outcome, an instance, or a data point. In general each x_j is a vector of values, and θ is a vector of real-valued parameters. For example, for a Gaussian distribution $\theta = \langle \mu, \sigma^2 \rangle$.

We assume that the examples are independent, so the probability of the set is the product of the probabilities of the individual examples:

$$f(x_1, \dots, x_n; \theta) = \prod_j f_\theta(x_j; \theta).$$

The notation above makes us think of the distribution θ as fixed and the examples x_j as unknown, or varying. However, we can think of the training data as fixed and consider alternative parameter values. This is the point of view behind the definition of the likelihood function:

$$L(\theta; x_1, \dots, x_n) = f(x_1, \dots, x_n; \theta).$$

Note that if $f(x; \theta)$ is a probability mass function, then the likelihood is always less than one, but if $f(x; \theta)$ is a probability density function, then the likelihood can be greater than one, since densities can be greater than one.

The principle of maximum likelihood says that given the training data, we should use as our model the distribution $f(\cdot; \hat{\theta})$ that gives the greatest possible probability to the training data. Formally,

$$\hat{\theta} = \operatorname{argmax}_{\theta} L(\theta; x_1, \dots, x_n).$$

The value $\hat{\theta}$ is called the maximum likelihood estimator (MLE) of θ . In general the hat notation indicates an estimated quantity; if necessary we will use notation like $\hat{\theta}_{\text{MLE}}$ to indicate the nature of an estimate.

2 Examples of maximizing likelihood

As a first example of finding a maximum likelihood estimator, consider estimating the parameter of a Bernoulli distribution. A random variable with this distribution is a formalization of a coin toss. The value of the random variable is 1 with probability θ and 0 with probability $1 - \theta$. Let X be a Bernoulli random variable, and let x be an outcome of X . We have

$$P(X = x) = \begin{cases} \theta & \text{if } x = 1 \\ 1 - \theta & \text{if } x = 0 \end{cases}$$

Usually, we use the notation $P(\cdot)$ for a probability mass, and the notation $p(\cdot)$ for a probability density. For mathematical convenience write $P(X)$ as

$$P(X = x) = \theta^x (1 - \theta)^{1-x}.$$

Suppose that the training data are x_1 through x_n where each $x_j \in \{0, 1\}$. The likelihood function is

$$L(\theta; x_1, \dots, x_n) = f(x_1, \dots, x_n; \theta) = \prod_{i=1}^n P(X = x_i) = \theta^h (1 - \theta)^{n-h}$$

where $h = \sum_i x_i$. The maximization is over the possible scalar values $0 \leq \theta \leq 1$.

We can do the maximization by setting the derivative with respect to θ equal to zero. The derivative is

$$\begin{aligned} \frac{d}{d\theta} \theta^h (1 - \theta)^{n-h} &= h\theta^{h-1} (1 - \theta)^{n-h} + \theta^h (n - h) (1 - \theta)^{n-h-1} (-1) \\ &= \theta^{h-1} (1 - \theta)^{n-h-1} [h(1 - \theta) - (n - h)\theta] \end{aligned}$$

which has solutions $\theta = 0$, $\theta = 1$, and $\theta = h/n$. The solution which is a maximum is clearly $\theta = h/n$ while $\theta = 0$ and $\theta = 1$ are minima. So we have the maximum likelihood estimate $\hat{\theta} = h/n$.

The log likelihood function, written $l(\cdot)$, is simply the logarithm of the likelihood function $L(\cdot)$. Because logarithm is a monotonic strictly increasing function, maximizing the log likelihood is precisely equivalent to maximizing the likelihood, or to minimizing the negative log likelihood.

For an example of minimizing the negative log likelihood (NLL), consider the problem of estimating the parameters of a univariate Gaussian distribution. This distribution is

$$f(x; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{(x - \mu)^2}{2\sigma^2}\right].$$

The NLL for one example x is

$$l(\mu, \sigma^2; x) = \log L(\mu, \sigma^2; x) = -\log \sigma - \log \sqrt{2\pi} - \frac{(x - \mu)^2}{2\sigma^2}.$$

Suppose that we have training data $\{x_1, \dots, x_n\}$. The maximum likelihood estimates are

$$\langle \hat{\mu}, \hat{\sigma}^2 \rangle = \operatorname{argmin}_{\langle \mu, \sigma^2 \rangle} \left[-n \log \sigma - n \log \sqrt{2\pi} - \frac{1}{2\sigma^2} \sum_i (x_i - \mu)^2 \right].$$

This expression is to be minimized simultaneously over two variables, but we can simplify it into two sequential univariate minimizations. The first is

$$\hat{\mu} = \operatorname{argmin}_{\mu} \sum_i (x_i - \mu)^2$$

while the second is

$$\hat{\sigma}^2 = \operatorname{argmin}_{\sigma^2} \left[-n \log \sigma \sqrt{2\pi} - \frac{1}{2\sigma^2} T \right]$$

where $T = \sum_i (x_i - \hat{\mu})^2$. In order to do the first minimization, write $(x_i - \mu)^2$ as $(x_i - \bar{x} + \bar{x} - \mu)^2$. Then

$$\sum_i (x_i - \mu)^2 = \sum_i (x_i - \bar{x})^2 - 2(\bar{x} - \mu) \sum_i (x_i - \bar{x}) + n(\bar{x} - \mu)^2.$$

The first term $\sum_i (x_i - \bar{x})^2$ does not depend on μ so it is irrelevant to the minimization. The second term equals zero, because $\sum_i (x_i - \bar{x}) = 0$. The third term is always positive, so it is clear that it is minimized when $\mu = \bar{x}$.

To perform the second minimization, work out the derivative symbolically and then work out when it equals zero:

$$\begin{aligned} \frac{\partial}{\partial \sigma} [-n \log \sigma - n \log \sqrt{2\pi} - \frac{1}{2} \sigma^{-2} T] &= -n\sigma^{-1} - \frac{1}{2}(-2\sigma^{-3})T \\ &= \sigma^{-1}(-n + T\sigma^{-2}) \\ &= 0 \text{ if } \sigma^2 = T/n. \end{aligned}$$

Maximum likelihood estimators are typically reasonable, but they may have issues. Consider the Gaussian variance estimator $\hat{\sigma}_{\text{MLE}}^2 = \sum_i (x_i - \bar{x})^2/n$ and the case where $n = 1$. In this case $\hat{\sigma}_{\text{MLE}}^2 = 0$. This estimate is guaranteed to be too small. Intuitively, the estimate is optimistically assuming that all future data points x_2 and so on will equal x_1 exactly.

It can be proved that in general the maximum likelihood estimate of the variance of a Gaussian is too small, on average:

$$E\left[\frac{1}{n} \sum_i (x_i - \bar{x})^2; \mu, \sigma^2\right] = \frac{n-1}{n} \sigma^2 < \sigma^2.$$

This phenomenon can be considered an instance of overfitting: the observed spread around the observed mean \bar{x} is less than the unknown true spread σ^2 around the unknown true mean μ .

3 Conditional likelihood

An important extension of the idea of likelihood is *conditional* likelihood. The conditional likelihood of θ given data x and y is $L(\theta; y|x) = f(y|x; \theta)$. Intuitively, y follows a probability distribution that is different for different x , but x itself is never unknown, so there is no need to have a probabilistic model of it. Technically, for each x there is a different distribution $f(y|x; \theta)$ of y , but all these distributions share the same parameters θ .

Given training data consisting of $\langle x_i, y_i \rangle$ pairs, the principle of maximum conditional likelihood says to choose a parameter estimate $\hat{\theta}$ that maximizes the product $\prod_i f(y_i|x_i; \theta)$. Note that we do not need to assume that the x_i are independent in order to justify the conditional likelihood being a product; we just need to assume that the y_i are independent when each is conditioned on its own x_i . For any specific value of x , $\hat{\theta}$ can then be used to predict values for y ; we assume that we never want to predict values of x .

Suppose that y is a binary (Bernoulli) outcome and that x is a real-valued vector. We can assume that the distribution of y is a fixed nonlinear function of a linear function of x . Specifically, we assume the conditional model

$$p(y|x; \alpha, \beta) = \frac{1}{1 + \exp -[\alpha + \sum_{j=1}^d \beta_j x_j]}.$$

This model is called logistic regression. We use j to index over the feature values x_1 to x_d of a single example of dimensionality d , since we use i below to index over training examples 1 to n . If necessary, the notation x_{ij} means the j th feature value of the i th example. Be sure to understand the distinction between a feature and a value of a feature. Essentially a feature is a random variable, while a value of a feature is a possible outcome of the random variable. Features may also be called attributes, predictors, or independent variables. The dependent random variable Y is sometimes called a dependent variable.

The logistic regression model is easier to understand in the form

$$\log \frac{p}{1-p} = \alpha + \sum_{j=1}^d \beta_j x_j$$

where p is an abbreviation for $p(y|x; \alpha, \beta)$. The ratio $p/(1-p)$ is called the odds of the event y given x , and $\log[p/(1-p)]$ is called the log odds. Since probabilities range between 0 and 1, odds range between 0 and $+\infty$ and log odds range unboundedly between $-\infty$ and $+\infty$. A linear expression of the form $\alpha + \sum_j \beta_j x_j$ can also take unbounded values, so it is reasonable to use a linear expression as a model for log odds, but not as a model for odds or for probabilities. Essentially, logistic regression is the simplest reasonable model for a random yes/no outcome that depends linearly on predictors x_1 to x_d .

For each feature j , $\exp(\beta_j x_j)$ is a multiplicative scaling factor on the odds $p/(1-p)$. If the predictor x_j is binary, then $\exp(\beta_j)$ is the extra odds of having the outcome $y = 1$ when $x_j = 1$, compared to when $x_j = 0$. If the predictor x_j is real-valued, then $\exp(\beta_j)$ is the extra odds of having the outcome $y = 1$ when the value of x_j increases by one unit. A major limitation of the basic logistic regression model is that the probability p must either increase monotonically, or decrease monotonically, as a function of each predictor x_j . The basic model does not allow the probability to depend in a U-shaped way on any x_j .

Given the training set $\{\langle x_1, y_1 \rangle, \dots, \langle x_n, y_n \rangle\}$, we learn a logistic regression classifier by maximizing the log joint conditional likelihood. This is the sum of

the log conditional likelihood for each training example:

$$LCL = \sum_{i=1}^n \log L(\theta; y_i|x_i) = \sum_{i=1}^n \log f(y_i|x_i; \theta).$$

Given a single training example $\langle x_i, y_i \rangle$, the log conditional likelihood is $\log p_i$ if the true label $y_i = 1$ and $\log(1 - p_i)$ if $y_i = 0$, where $p_i = p(y = 1|x_i; \theta)$.

To simplify the following discussion, assume from now on that $\alpha = \beta_0$ and $x_0 = 1$ for every example x , so the parameter vector θ is $\beta \in \mathbb{R}^{d+1}$. By grouping together the positive and negative training examples, we can write the total conditional log likelihood as

$$LCL = \sum_{i:y_i=1} \log p_i + \sum_{i:y_i=0} \log(1 - p_i).$$

The partial derivative of LCL with respect to parameter β_j is

$$\sum_{i:y_i=1} \frac{\partial}{\partial \beta_j} \log p_i + \sum_{i:y_i=0} \frac{\partial}{\partial \beta_j} \log(1 - p_i).$$

For an individual training example $\langle x, y \rangle$, if its label $y = 1$ the partial derivative is

$$\frac{\partial}{\partial \beta_j} \log p = \frac{1}{p} \frac{\partial}{\partial \beta_j} p$$

while if $y = 0$ it is

$$\frac{\partial}{\partial \beta_j} \log(1 - p) = \frac{1}{1 - p} \left(- \frac{\partial}{\partial \beta_j} p \right).$$

Let $e = \exp[-\sum_{j=0}^d \beta_j x_j]$ so $p = 1/(1 + e)$ and $1 - p = (1 + e - 1)/(1 + e) = e/(1 + e)$. With this notation we have

$$\begin{aligned} \frac{\partial}{\partial \beta_j} p &= (-1)(1 + e)^{-2} \frac{\partial}{\partial \beta_j} e \\ &= (-1)(1 + e)^{-2} (e) \frac{\partial}{\partial \beta_j} [-\sum_j \beta_j x_j] \\ &= (-1)(1 + e)^{-2} (e)(-x_j) \\ &= \frac{1}{1 + e} \frac{e}{1 + e} x_j \\ &= p(1 - p)x_j. \end{aligned}$$

So

$$\frac{\partial}{\partial \beta_j} \log p = (1 - p)x_j \quad \text{and} \quad \frac{\partial}{\partial \beta_j} \log(1 - p) = -px_j.$$

For the entire training set the partial derivative of the log conditional likelihood with respect to β_j is

$$\frac{\partial}{\partial \beta_j} LCL = \sum_{i:y_i=1} (1 - p_i)x_{ij} + \sum_{i:y_i=0} -p_i x_{ij} = \sum_i (y_i - p_i)x_{ij}$$

where x_{ij} is the value of the j th feature of the i th training example. Setting the partial derivative to zero yields

$$\sum_i y_i x_{ij} = \sum_i p_i x_{ij}.$$

We have one equation of this type for each parameter β_j . The equations can be used to check the correctness of a trained model.

Informally, but not precisely, the expression $\sum_i y_i x_{ij}$ is the average value over the training set of the i th feature, where each training example is weighted 1 if its true label is positive, and 0 otherwise. The expression $\sum_i p_i x_{ij}$ is the same average, except that each example i is weighted according to its predicted probability p_i of being positive. When the logistic regression classifier is trained correctly, then these two averages must be the same for every feature. The special case for $j = 0$ gives

$$\frac{1}{n} \sum_i y_i = \frac{1}{n} \sum_i p_i.$$

In words, the empirical base rate probability of being positive must equal the average predicted probability of being positive.

4 Stochastic gradient training

There are several sophisticated ways of actually doing the maximization of the total conditional log likelihood, i.e. the conditional log likelihood summed over all training examples $\langle x_i, y_i \rangle$; for details see [Minka, 2007, Komarek and Moore, 2005]. However, here we consider a method called stochastic gradient ascent. This method changes the parameter values to increase the log likelihood based on one

example at a time. It is called stochastic because the derivative based on a randomly chosen single example is a random approximation to the true derivative based on all the training data.

Consider a single training example $\langle x, y \rangle$, where again we drop the subscript i for convenience. Consider the j th parameter for $0 \leq j \leq d$. The partial derivative of the log likelihood given this single example is

$$\frac{\partial}{\partial \beta_j} \log L(\beta; x, y) = (y - p)x_j$$

where $y = 1$ or $y = 0$. For each j , we increase the log likelihood incrementally by doing the update

$$\beta_j := \beta_j + \lambda(y - p)x_j.$$

Here λ is a multiplier called the learning rate that controls the magnitude of the changes to the parameters.

Stochastic gradient ascent (or descent, for a minimization problem) is a method that is often useful in machine learning. Experience suggests some heuristics for making it work well in practice.

- The training examples are sorted in random order, and the parameters are updated for each example sequentially. One complete update for every example is called an epoch. Typically, a small constant number of epochs is used, perhaps 3 to 100 epochs.
- The learning rate is chosen by trial and error. It can be kept constant across all epochs, e.g. $\lambda = 0.1$ or $\lambda = 1$, or it can be decreased gradually as a function of the epoch number.
- Because the learning rate is the same for every parameter, it is useful to scale the features x_j so that their magnitudes are similar for all j . Given that the feature x_0 has constant value 1, it is reasonable to normalize every other feature to have mean zero and variance 1, for example.

Stochastic gradient ascent (or descent) has some properties that are very useful in practice. First, suppose that $x_j = 0$ for most features j of a training example x . Then updating β_j based on x can be skipped. This means that the time to do one epoch is $O(nfd)$ where n is the number of training examples, d is the number of features, and f is the average number of nonzero feature values per example. If an example x is the bag-of-words representation of document, then d is the size of

the vocabulary (often over 30,000) but fd is the average number of words actually used in a document (often under 300).

Second, suppose that the number n of training examples is very large, as is the case in many modern applications. Then, a stochastic gradient method may converge to good parameter estimates in less than one epoch of training. In contrast, a training method that computes the log likelihood of all data and uses this in the same way regardless of n will be inefficient in how it uses the data.

For each example, a stochastic gradient method updates all parameters once. The dual idea is to update one parameter at a time, based on all examples. This method is called coordinate ascent (or descent). For feature j the update rule is

$$\beta_j := \beta_j + \lambda \sum_i (y_i - p_i) x_{ij}.$$

The update for the whole parameter vector $\bar{\beta}$ is

$$\bar{\beta} := \bar{\beta} + \lambda(\bar{y} - \bar{p})^T X$$

where the matrix X is the entire training set and the column vector \bar{y} consists of the 0/1 labels for every training example. Often, coordinate ascent converges too slowly to be useful. However, it can be useful to do one update of $\bar{\beta}$ after all epochs of stochastic gradient ascent.

Regardless of the method used to train a model, it is important to remember that optimizing the model perfectly on the training data usually does not lead to the best possible performance on test examples. There are several reasons for this:

- The model with best possible performance may not belong to the family of models under consideration. This is an instance of the principle “you cannot learn it if you cannot represent it.”
- The training data may not be representative of the test data, i.e. the training and test data may be samples from different populations.
- Fitting the training data as closely as possible may simply be overfitting.
- The objective function for training, namely log likelihood or conditional log likelihood, may not be the desired objective from an application perspective; for example, the desired objective may be classification accuracy.

5 Regularization

Consider learning a logistic regression classifier for categorizing documents. Suppose that word number j appears only in documents whose labels are positive. The partial derivative of the log conditional likelihood with respect to the parameter for this word is

$$\frac{\partial}{\partial \beta_j} LCL = \sum_i (y_i - p_i) x_{ij}.$$

This derivative will always be positive, as long as the predicted probability p_i is not perfectly one for all these documents. Therefore, following the gradient will send β_j to infinity. Then, every test document containing this word will be predicted to be positive with certainty, regardless of all other words in the test document. This over-confident generalization is an example of overfitting.

There is a standard method for solving this overfitting problem that is quite simple, but quite successful. The solution is called regularization. The idea is to impose a penalty on the magnitude of the parameter values. This penalty should be minimized, in a trade-off with maximizing likelihood. Mathematically, the optimization problem to be solved is

$$\hat{\beta} = \operatorname{argmax}_{\beta} LCL - \mu \|\beta\|_2$$

where $\|\beta\|_2$ is the L_2 norm of the parameter vector, and the constant μ quantifies the trade-off between maximizing likelihood and minimizing parameter values.

This type of regularization is called quadratic or Tikhonov regularization. A major reason why it is popular is that it can be derived from several different points of view. In particular, it arises as a consequence of assuming a Gaussian prior on parameters. It also arises from theorems on minimizing generalization error, i.e. error on independent test sets drawn from the same distribution as the training set. And, it arises from robust classification: assuming that each training point lies in an unknown location inside a sphere centered on its measured location.

Stochastic gradient following is easily extended to include regularization. We simply include the penalty term when calculating the gradient for each example. Consider

$$\frac{\partial}{\partial \beta_j} [\log p(y|x; \beta) - \mu \sum_{j=0}^d \beta_j^2] = \left[\frac{\partial}{\partial \beta_j} \log p(y|x; \beta) \right] - \mu 2\beta_j.$$

Remember that for logistic regression the partial derivative of the log conditional

likelihood for one example is

$$\frac{\partial}{\partial \beta_j} \log p(y|x; \beta) = (y - p)x_j$$

so the update rule with regularization is

$$\beta_j := \beta_j + \lambda[(y - p)x_j - 2\mu\beta_j]$$

where λ is the learning rate. Update rules like the one above are often called “weight decay” rules, since the weight β_j is made smaller at each update unless $y - p$ has the same sign as x_j .

Straightforward stochastic gradient ascent for training a regularized logistic regression model loses the desirable sparsity property described above, because the value of every parameter β_j must be decayed for every training example. How to overcome this computational inefficiency is described in [Carpenter, 2008].

Writing the regularized optimization problem as a minimization gives

$$\hat{\beta} = \operatorname{argmin}_{\beta} \sum_{i=1}^n -\log p(y_i|x_i; \beta) + \mu \sum_{j=0}^d \beta_j^2.$$

The expression $-\log p(y_i|x_i; \beta)$ is called the “loss” for training example i . If the predicted probability, using β , of the true label y_i is close to 1, then the loss is small. But if the predicted probability of y_i is close to 0, then the loss is large. Losses are always non-negative; we want to minimize them. We also want to minimize the numerical magnitude of the trained parameters.

References

- [Carpenter, 2008] Carpenter, B. (2008). Lazy sparse stochastic gradient descent for regularized multinomial logistic regression. Technical report, Alias-i. Available at <http://lingpipe-blog.com/lingpipe-white-papers/>.
- [Komarek and Moore, 2005] Komarek, P. and Moore, A. W. (2005). Making logistic regression a core data mining tool with TR-IRLS. In *Proceedings of the Fifth IEEE International Conference on Data Mining (ICDM'05)*, pages 685–688.

[Minka, 2007] Minka, T. P. (2007). A comparison of numerical optimizers for logistic regression. First version 2001. Unpublished paper available at <http://research.microsoft.com/~minka>.

CSE 250B Quiz 3, January 21, 2010

Write your name:

Let $f_\theta(x; \theta)$ where $x \in \mathbb{R}$ be the probability density function (pdf) of the uniform distribution over the range 0 to θ . Precisely, $f_\theta(x; \theta) = 1/\theta$ if $0 \leq x \leq \theta$ while $f_\theta(x; \theta) = 0$ otherwise.

Let x_1 to x_n be an independent identically distributed (iid) sample from f_θ for some unknown true parameter value $\theta > 0$. The maximum likelihood estimator (MLE) of θ is $\hat{\theta} = \max_i x_i$.

[3 points] In one or two sentences, explain intuitively the reason why this is the MLE. You do not need to use any equations.

Note: The MLE above is an example of overfitting, since the true value of θ is almost certainly larger than the MLE.

CSE 250B Quiz 4, January 27, 2011

Write your name:

Assume that winning or losing a basketball game is similar to flipping a biased coin. Suppose that San Diego State University (SDSU) has won all six games that it has played.

- (a) The maximum likelihood estimate (MLE) of the probability that SDSU will win its next game is 1.0. Explain why this is the MLE. (Using equations is not required.)
- (b) This MLE can be viewed as overfitting. Explain why.

CSE 250B Quiz 4, January 28, 2010

The objective function to be minimized when training an L_2 -regularized linear regression model is

$$E = \sum_{i=1}^n (f(x_i; w) - y_i)^2 + \mu \sum_{j=0}^d w_j^2$$

where the model is

$$f(x_i; w) = \sum_{j=0}^d w_j x_{ij}.$$

All notation above is the same as in the class lecture notes.

[3 points] Work out the partial derivative of the objective function with respect to weight w_j .

Answer. Let f_i be an abbreviation for $f(x_i; w)$ and let d_i be an abbreviation for $\frac{\partial}{\partial w_j} f_i$. Note that $d_i = x_{ij}$. The expanded objective function is

$$E = \left[\sum_{i=1}^n f_i^2 - 2f_i y_i + y_i^2 \right] + \mu \sum_{j=0}^d w_j^2.$$

The partial derivative is

$$\frac{\partial}{\partial w_j} E = \left[\sum_{i=1}^n [2f_i d_i - 2y_i d_i + 0] \right] + 2\mu w_j$$

which is

$$\frac{\partial}{\partial w_j} E = 2[\mu w_j + \sum_{i=1}^n (f_i - y_i) x_{ij}].$$

Additional note. Because the goal is to minimize E , we do gradient descent, not ascent, with the update rule

$$w_j := w_j - \lambda \frac{\partial}{\partial w_j} E.$$

The update rule says that if the average over all training examples i of $(f_i - y_i)x_{ij}$ is positive, then w_j must be decreased. Assume that x_{ij} is non-negative; this update rule is reasonable because then f_i is too big on average, and decreasing w_j will make f_i decrease. The update rule also says that, because of regularization, w_j must always be decreased even more, by $2\lambda\mu$ times its current value.

CSE 250B Quiz 4, January 27, 2011

Write your name:

Assume that winning or losing a basketball game is similar to flipping a biased coin. Suppose that San Diego State University (SDSU) has won all six games that it has played.

- (a) The maximum likelihood estimate (MLE) of the probability that SDSU will win its next game is 1.0. Explain why this is the MLE. (Using equations is not required.)
- (b) This MLE can be viewed as overfitting. Explain why.

CSE 250B Quiz, February 3, 2011

Your name:

Suppose you are doing stochastic gradient descent to minimize the following error function on a single training example:

$$E = e(f(x; w) - y)$$

Work out the stochastic gradient update rule as specifically as possible, when the error function is absolute error: $e(z) = |z|$.

Hint: Use the notation $e'(z)$ for the derivative of $e(z)$.

Answer: For each component w_j of the parameter vector w , the update rule is

$$w_j := w_j - \alpha \frac{\partial}{\partial w_j} E$$

where α is the learning rate. We can work out

$$\frac{\partial}{\partial w_j} E = e'(f(x; w) - y) \frac{\partial}{\partial w_j} f$$

If $e(z) = |z|$ then $e'(z) = \text{sign}(z)$ so

$$w_j := w_j - \alpha \text{sign}(f(x; w) - y) \frac{\partial}{\partial w_j} f.$$

Intuitively, if $f(x; w)$ is too large, and increasing w_j makes f increase, then w_j should be decreased.

Project assignment 2

For this assignment you may work in a team of two with one partner, or in a team of three. For multiple reasons, working alone or in a larger team is not appropriate. The joint report for your team must be submitted in hard copy at the start of class on Thursday, February 3, 2011.

This project uses data published by Kevin Hillstrom, a well-known data mining consultant. You can find the data at <http://cseweb.ucsd.edu/users/elkan/250B/HillstromData.csv>. A paper by Nicholas Radcliffe available at <http://www.stochasticolutions.com/etailPaper.html> is a good previous analysis of this dataset. Sections 1 through 5.2 are readable and useful, while Sections 5.3 and later are harder to follow and less relevant.

You should build three separate models: one for customers who are not sent any email promotion, one for customers who are sent the “men’s clothing” email, and one for customers who are sent the “women’s clothing” email. Each model consists of two regularized logistic regression classifiers and one regression function. The three submodels are part of the overall model

$$E[\text{spend}|x, \text{treatment}] = E[\text{spend}|\text{purchase}, x, \text{treatment}] \cdot p(\text{purchase}|\text{visit}, x, \text{treatment}) \cdot p(\text{visit}|x, \text{treatment}).$$

In your report, explain the equation above carefully. Note that

- x is the vector of attribute values describing a customer,
- “women’s clothing” email, “men’s clothing” email, and no email are the three possible values of the random variable “treatment,”
- “visit” and “purchase” are binary random variables that have a certain probability of being true or false for each customer, and
- “spend” is a real-valued random variable for each customer.

You should implement your own Matlab code for training logistic regression with regularization. The optimization method should be stochastic gradient descent (SGD). Your implementation should handle discrete features as well as real-valued features, and up to 30,000 training examples with up to 20 features. One difficulty with SGD is choosing the learning rate λ . You can find advice at <http://leon.bottou.org/projects/sgd>.

For the real-valued regression part of your model, you may use regularized linear regression or any standard Matlab function such as `nlinfit`. Anywhere

that you apply regularization, use cross-validation to select a good value for its strength. In your report, analyze experimentally and theoretically the big-O time and space complexity of your implementation. Explain any non-trivial obstacles that you did or did not overcome in trying to achieve good time complexity.

You must think carefully about how to do experiments that are conceptually sound. Your report should contain three separate sections that explain the design of your experiments, the results of the experiments, and conclusions drawn from the experiments. In the concluding section, provide a well-founded quantitative estimate of the benefit achievable by sending each different email to an optimal subset of future customers. Also estimate the benefit of doing data mining compared to a no-learning strategy. Make the concluding section understandable and useful for a manager who is not interested in technical questions.

You may want to use a scripting language to turn the raw data into a numerical dataset for input into Matlab. Describe any preprocessing and feature selection that you do in your report.