

Linear Regression

Sunday, August 9, 2020 6:01 PM

Data:

	Lot.Area	Bedroom.AbvGr	SalePrice
0	31770	3	215.0
1	11622	2	105.0
2	14267	3	172.0
3	11160	3	244.0
4	13830	3	189.9
...
96	3182	2	151.0
97	2544	2	149.5
98	2544	2	152.0
99	4403	2	222.0
100	2117	3	177.5

101 rows × 3 columns

Notation : Lot Area = α_1 , Bedroom.AbvGr = α_2

Sale price = y

Training set : $(\alpha^{(i)}, y^{(i)})$ where $i = 0 : 100$
 $i \in (X, Y)$

Hypothesis : $h(\alpha) = \theta_0 + \theta_1 \alpha_1 + \theta_2 \alpha_2$

$$= \theta_0 \alpha_0 + \theta_1 \alpha_1 + \theta_2 \alpha_2$$

$$= \sum_{i=0}^n \theta_i \alpha_i = \theta^T \alpha$$

$n = 100$ of training examples

① Iterative : $\theta = \theta + \alpha \sum_{i=1}^n (y^{(i)} - h_\theta(\alpha^{(i)})) \alpha^{(i)}$

$$\text{Since, } \frac{\partial}{\partial \theta} J(\theta) = \left(h_{\theta}(x^{(i)}) - y^{(i)} \right) x^{(i)}$$

where, $J(\theta)$ is mean squared error
 (cost function / loss)

② Implementation of Normal Equation (closed form)

$$X = \begin{bmatrix} x_0 & x_1 & x_2 \\ 1 & * & * \\ \vdots & * & * \\ \vdots & \vdots & \vdots \\ 1 & 1 & : \end{bmatrix} \quad \begin{array}{l} \text{~n examples} \\ \sim \end{array} \quad \begin{bmatrix} -(x^{(1)})^T \\ -(x^{(2)})^T \\ \vdots \\ -(x^{(n)})^T \end{bmatrix}$$

$$Y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

$$h_{\theta}(x^{(i)}) = X\theta - Y$$

$$J(\theta) = \frac{1}{2} (X\theta - Y)^T (X\theta - Y)$$

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \frac{1}{2} (X\theta - Y)^T (X\theta - Y)$$

$$= \frac{1}{2} \nabla_{\theta} (X\theta)^T X\theta - (X\theta)^T Y \\ - Y^T (X\theta) + Y^T Y$$

$$= \frac{1}{2} \nabla_{\theta} (\theta^T (x^T x) \theta - 2 (x^T y)^T \theta)$$

$$= \frac{1}{2} (2 x^T x \theta - 2 x^T y)$$

$$= x^T x \theta - x^T y$$

To find min arg θ ,

$$\nabla_{\theta} J(\theta) = 0 \Rightarrow x^T x \theta - x^T y = 0$$

$$\theta = (x^T x)^{-1} x^T y$$

[Using, $a^T b = b^T a$, $\nabla_{\alpha} b^T a = b$ $\nabla_{\alpha} x^T A \alpha = 2 A \alpha$
for a symmetric matrix 'A']

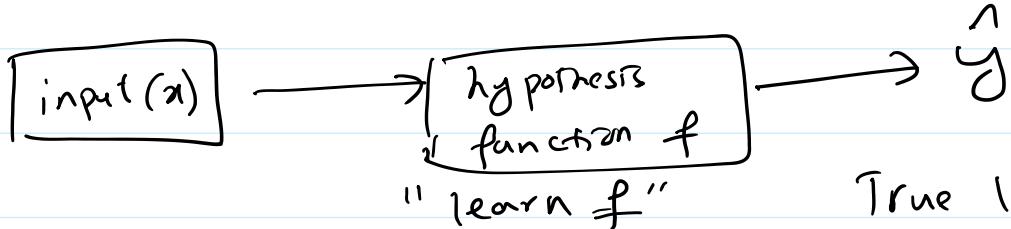
Code Implementation Note

Define X and Y, then use normal equation to find ' θ ' parameter directly.

LogisticRegression

Tuesday, August 11, 2020 9:08 AM

(A) Derivation of logistic Regression



True label $\rightarrow y$

$$\text{Div}(\hat{y}, y)$$

is loss.

$$y|x; \theta \sim \text{Bernoulli}(\phi)$$

or generally \sim Exponential family (η)

$$\text{PMF} \quad p(y|n) = b(y) \exp(n^T T(y) - a(n))$$

$$h(x) = E[y|x]$$

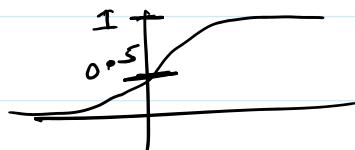
$$\text{on } h_{\theta}(x) = E[y|x; \theta] = \phi = \frac{1}{1 + e^{-\theta^T x}}$$

$$= \frac{1}{1 + e^{-\theta^T x}}$$

$g(\eta) = E[T(y); \eta]$ is canonical response function

In our case Sigmoid.

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}} \quad \text{when } g \text{ is sigmoid}$$



PMF

$$p(y|x; \theta) = \phi^y (1-\phi)^{1-y}$$

$$= (h_{\theta}(x))^y (1 - h_{\theta}(x))^{1-y}$$

Likelihood (function of θ)

$$L(\theta) = P(y|x; \theta)$$

$$= \prod_{i=1}^n P(y^{(i)}|x^{(i)}; \theta)$$

$y^{(i)} = \theta^\top x^{(i)} + \epsilon^{(i)}$
 assumption,
 $\epsilon^{(i)}$ are IID
 distributed

$$= \prod_{i=1}^n (h_{\theta}(x^{(i)}))^{y^{(i)}} (1 - h_{\theta}(x^{(i)}))^{1-y^{(i)}}$$

Since \log is strictly increasing function,
 we could find θ that maximize $\log(L(\theta))$

So,

$$\ell(\theta) = \log L(\theta) = \sum_{i=1}^n y^{(i)} \log h_{\theta}(x^{(i)}) + (-y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))$$

$$\frac{\partial}{\partial \theta_j} \ell(\theta) = (y - h_{\theta}(x)) x_j$$

$$\text{Note } [g'(x) = g(x)(1 - g(x))]$$

$$\text{So, } \theta := \theta + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) x$$

③ Problem set 1(A) and 1(B)

(a) [10 points] In lecture we saw the average empirical loss for logistic regression:

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h_\theta(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)})),$$

where $y^{(i)} \in \{0, 1\}$, $h_\theta(x) = g(\theta^T x)$ and $g(z) = 1/(1 + e^{-z})$.

Find the Hessian H of this function, and show that for any vector z , it holds true that

$$z^T Hz \geq 0.$$

Hint: You may want to start by showing that $\sum_i \sum_j z_i x_i x_j z_j = (x^T z)^2 \geq 0$. Recall also that $g'(z) = g(z)(1 - g(z))$.

Remark: This is one of the standard ways of showing that the matrix H is positive semi-definite, written " $H \succeq 0$." This implies that J is convex, and has no local minima other than the global one. If you have some other way of showing $H \succeq 0$, you're also welcome to use your method instead of the one above.

Answer:

$$\nabla_\theta J(\theta) = \begin{bmatrix} & & \\ & & \\ \vdots & & J(\theta) \\ & & \end{bmatrix}$$

For simplicity assume we have only one training example
on the set so we have

$$\begin{aligned} \frac{\partial J(\theta)}{\partial \theta} &= \frac{\partial}{\partial \theta} (y \log(h_\theta(x)) + (1-y) \log(1-h_\theta(x))) \\ &= \frac{y}{h_\theta(x)} \frac{\partial}{\partial \theta} h_\theta(x) + (1-y) \frac{\partial}{\partial \theta} (1-h_\theta(x)) \\ &= g'(\theta^T x) x^T \left(\frac{y}{h_\theta(x)} + \frac{1-y}{1-h_\theta(x)} \right) \\ &= g'(\theta^T x) x^T \left(\frac{y-h_\theta(x)}{h_\theta(x)(1-h_\theta(x))} \right) \\ &= g'(\theta^T x) x^T \frac{(y-h_\theta(x))}{g(\theta^T x)} \\ &= (y-h_\theta(x)) x^T \end{aligned}$$

\Rightarrow st $\frac{\partial J(\theta)}{\partial \theta} = (y-h_\theta(x)) x^T$

Now

$$H_{ij} = x_j h_\theta(x) (h_\theta(x) + 1) x_i$$

$$Hz = \sum_j H_{ij} x_j x_i = \sum_j h_\theta(x) (h_\theta(x) + 1) x_i x_j$$

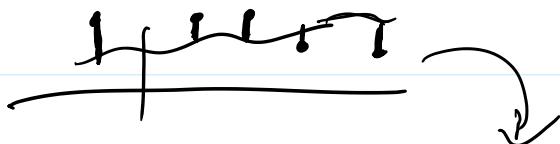
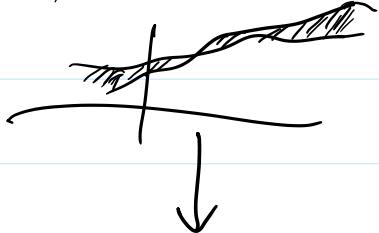
$$= h(n)(h(n)-1) \sum_{i,j} z_i z_j$$

$$= h(n)(h(n)-1) (\bar{z}^2)$$

≥ 0

So $\sum H \geq 0 \Rightarrow$ here, $H \geq 0$ so this P.D.O.S/H is Convex

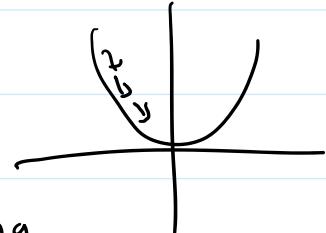
Use input sample to estimate error.



$$E[\text{div}(f(x; \theta), g(x))] \approx \frac{1}{n} \sum_{i=1}^n \text{div}(f(x_i; \theta), \hat{y}_i)$$

Implementation Notes (C)

Instead of iteratively going down hill by ' α ' learning rate (gradient descent),



We used quadratic approximation using Newton's method. For vector input x :

$$\theta := \theta - H^{-1} \nabla_\theta l(\theta) \quad (\text{Newton-Raphson method})$$

$$H \text{ is Hessian}, H_{ij} = \frac{\partial^2}{\partial \theta_i \partial \theta_j} l(\theta)$$

'numerical analysis'

Algorithm

Training

- ① Initialize the parameters θ

① Define $h_{\theta}(x) = g(\theta^T x)$ where g is sigmoid
 $y \in \{0, 1\}$ (in ps'1)

② Compute $\nabla J(\theta)$

empirical loss, $\frac{1}{n} \sum (y - h_{\theta}(x)) x$

③ Compute Hessian H ,

average) $\frac{1}{n} \sum h_{\theta}(x)(1-h_{\theta}(x)) x^T x$

④ Use Newton-Raphson, Repeat until all
 Convergence

Testing

⑤ To predict, feed the input x to
 hypothesis, $h_{\theta}(x)$ to produce
 prediction.

- (b) [5 points] **Coding problem.** Follow the instructions in `src/p01b.logreg.py` to train a logistic regression classifier using Newton's Method. Starting with $\theta = \vec{0}$, run Newton's Method until the updates to θ are small: Specifically, train until the first iteration k such that $\|\theta_k - \theta_{k-1}\|_1 < \epsilon$, where $\epsilon = 1 \times 10^{-5}$. Make sure to write your model's predictions to the file specified in the code.

Answer:

(A) Derivation of GDA model

Comparison with Discriminative algorithm:

In Disc we have to classify $y=0$ or $y=1$ (binary example) given features x . Logistic Reg tried to find a straight line (Decision boundary) that separates the two levels. While, GDA looks at how $y=1$ and $y=0$ are made up of and based the prediction of a example by comparing if it fits $y=1$ or $y=0$ properties.

Disc Algo: Find $P(y|x)$ directly

GDA : Find $P(x|y)$ and $P(y)$ Then use makes strong assumption from to to find $P(y|x)$ e.g., $P(x|y=0)$, $P(x|y=1)$
 $P(y) \rightarrow$ class prior

$$\text{Use bayes rule, } P(y|x) = \frac{P(x|y) P(y)}{P(x)}$$

$$\text{arg max}_y P(y|x) = \text{arg max}_y P(x|y) P(y)$$

Model:

$y \sim \text{Binomial}(\delta)$

$x|y=0 \sim N(\mu_0, \Sigma)$

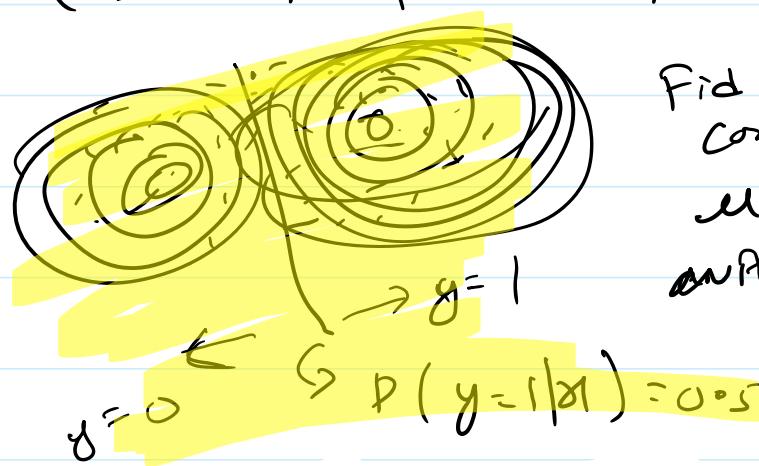
$x|y=1 \sim N(\mu_1, \Sigma)$

$P(y)$
 $P(x|y=0)$ }
 $P(x|y=1)$ } → find prob mass/density funct

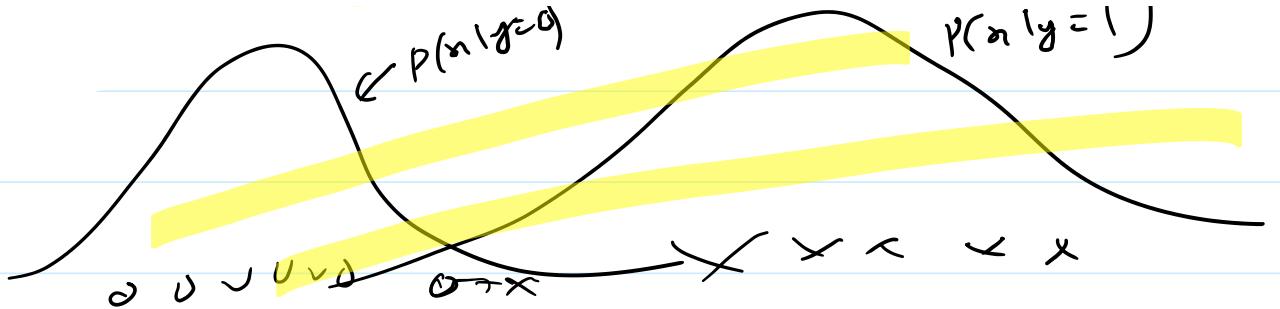
$$l(\phi, \mu_0, \mu_1, \Sigma) = \log \prod_{i=1}^n p(x^{(i)}, y^{(i)}; \phi, \mu_0, \mu_1, \Sigma)$$

$$= \log \prod_{i=1}^n \frac{p(x^{(i)}|y^{(i)}; \mu_0, \mu_1, \Sigma)}{p(y^{(i)}; \phi)}$$

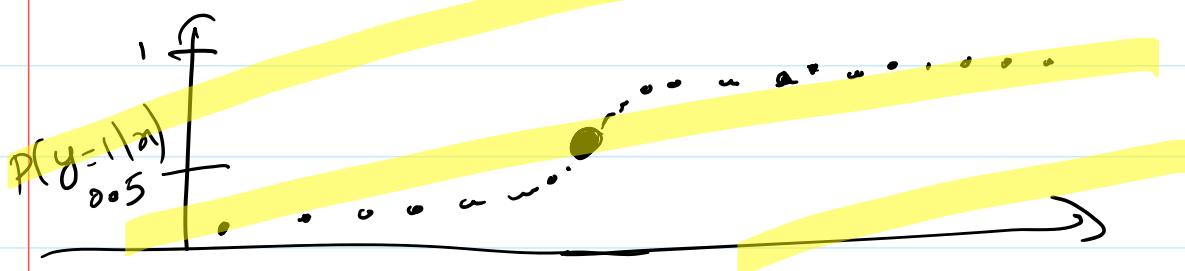
Take derivative of l wrt parameters ($\phi, \mu_0, \mu_1, \Sigma$) to find the parameter that maximize the likelihood. (Done in Problem Set below)



Fid gaussian controlled by
 μ_0, μ_1 , to shift
on the shape shape
 Σ .



$P(y=1) = 0.5$ since data is half and half



Generations

ym Binomial (\$)

$$x_{(j=0)} \sim N(\bar{x}_0, \Sigma) \Rightarrow$$

$$x|y=1 \sim N(\mu_1, \Sigma)$$

↳ complaints

$$P(y=1|x) = \frac{1}{1+e^{-\theta^T x}}$$

(B) Problem Set 1 (C) (D) (E)

- (c) [5 points] Recall that in GDA we model the joint distribution of (x, y) by the following equations:

$$p(y) = \begin{cases} \phi & \text{if } y = 1 \\ 1 - \phi & \text{if } y = 0 \end{cases}$$

$$p(x|y=0) = \frac{1}{(2\pi)^{n/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_0)^T \Sigma^{-1} (x - \mu_0)\right)$$

$$p(x|y=1) = \frac{1}{(2\pi)^{n/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1} (x - \mu_1)\right),$$

where ϕ, μ_0, μ_1 , and Σ are the parameters of our model.

Suppose we have already fit ϕ, μ_0, μ_1 , and Σ , and now want to predict y given a new point x . To show that GDA results in a classifier that has a linear decision boundary, show the posterior distribution can be written as

$$p(y=1 | x; \phi, \mu_0, \mu_1, \Sigma) = \frac{1}{1 + \exp(-(\theta^T x + \theta_0))},$$

where $\theta \in \mathbb{R}^n$ and $\theta_0 \in \mathbb{R}$ are appropriate functions of ϕ, Σ, μ_0 , and μ_1 .

Answer:

$$\begin{aligned} P(y=1 | x) &= \frac{P(n|y=1) P(y=1)}{P(n|y=1) P(y=1) + P(n|y=0) P(y=0)} \\ &= \frac{1}{1 + \frac{P(n|y=0) P(y=0)}{P(n|y=1) P(y=1)}} \\ &= \frac{1}{1 + \frac{\exp(-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1} (x - \mu_1)) / (\phi)}{\exp(-\frac{1}{2}(x - \mu_0)^T \Sigma^{-1} (x - \mu_0)) / (1 - \phi)}} \quad \rightarrow x \\ x &= e^{\left(-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1} (x - \mu_1) - \left(-\frac{1}{2}(x - \mu_0)^T \Sigma^{-1} (x - \mu_0)\right)\right) \log\left(\frac{1-\phi}{\phi}\right)} \\ x &= e^{\left(-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1} (x - \mu_1) + \frac{1}{2}(x - \mu_0)^T \Sigma^{-1} (x - \mu_0)\right) \log\left(\frac{1-\phi}{\phi}\right)} \end{aligned}$$

$$= \frac{1}{1 + \exp(-[(\mu_1 - \mu_0)^T \Sigma^{-1} x] + \frac{1}{2}(\mu_0^T \Sigma^{-1} \mu_0 - \mu_1^T \Sigma^{-1} \mu_1) - \log \frac{1-\phi}{\phi})}$$

So we can take $\theta = (\mu_1 - \mu_0)^T \Sigma^{-1} x$ $\phi = \frac{1-\phi}{\phi}$

$\Rightarrow e$
for coding part

Using $(A+B)C = (AC) + (BC)$ $A^T + B^T = (A+B)^T$

- (d) [7 points] For this part of the problem only, you may assume n (the dimension of x) is 1, so that $\Sigma = [\sigma^2]$ is just a real number, and likewise the determinant of Σ is given by $|\Sigma| = \sigma^2$. Given the dataset, we claim that the maximum likelihood estimates of the parameters are given by

$$\begin{aligned}\phi &= \frac{1}{m} \sum_{i=1}^m 1\{y^{(i)} = 1\} \\ \mu_0 &= \frac{\sum_{i=1}^m 1\{y^{(i)} = 0\}x^{(i)}}{\sum_{i=1}^m 1\{y^{(i)} = 0\}} \\ \mu_1 &= \frac{\sum_{i=1}^m 1\{y^{(i)} = 1\}x^{(i)}}{\sum_{i=1}^m 1\{y^{(i)} = 1\}} \\ \Sigma &= \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T\end{aligned}$$

The log-likelihood of the data is

$$\begin{aligned}\ell(\phi, \mu_0, \mu_1, \Sigma) &= \log \prod_{i=1}^m p(x^{(i)}, y^{(i)}; \phi, \mu_0, \mu_1, \Sigma) \\ &= \log \prod_{i=1}^m p(x^{(i)} | y^{(i)}; \mu_0, \mu_1, \Sigma) p(y^{(i)}; \phi).\end{aligned}$$

By maximizing ℓ with respect to the four parameters, prove that the maximum likelihood estimates of ϕ , μ_0 , μ_1 , and Σ are indeed as given in the formulas above. (You may assume that there is at least one positive and one negative example, so that the denominators in the definitions of μ_0 and μ_1 above are non-zero.)

Answer:

$$\begin{aligned}\ell &= \log \prod_{i=1}^m \left[P(y^{(i)}=0 | x^{(i)}, \phi, \mu_0, \mu_1, \Sigma) + P(y^{(i)}=1 | x^{(i)}, \phi, \mu_0, \mu_1, \Sigma) \right] \\ &= \log \prod_{i=1}^m \left[1\{y^{(i)}=0\} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x^{(i)} - \mu_0)^2}{2\sigma^2}\right) (1-\phi) + \dots \right] \\ &= \sum_{i=1}^m \log \left[1\{y^{(i)}=0\} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x^{(i)} - \mu_0)^2}{2\sigma^2}\right) (1-\phi) + \dots \right] \\ &= \sum_{i=1}^m \log \left[1\{y^{(i)}=0\} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x^{(i)} - \mu_0)^2}{2\sigma^2}\right) + 1\{y^{(i)}=1\} \log \left(\frac{1}{\sqrt{2\pi}\sigma} \right) \right]\end{aligned}$$

$$\frac{\partial \ell}{\partial \phi} = \sum_{i=1}^m \frac{\partial}{\partial \phi} \left[y^{(i)} \log(\phi) + (1-y^{(i)}) \log(1-\phi) \right] = 0$$

$$\sum_{i=1}^m \frac{y^{(i)}}{\phi} + \frac{1-y^{(i)}}{1-\phi} = 0$$

$$\sum_{i=1}^m y^{(i)} - \phi y^{(i)} + \phi - \phi y^{(i)} = 0$$

$$\sum_{i=1}^m y^{(i)} + \phi \sum_{i=1}^m (1-y^{(i)}) = 0 \Rightarrow \phi = \frac{\sum_{i=1}^m y^{(i)}}{\sum_{i=1}^m (1-y^{(i)})}$$

$$\sum_{i=1}^m y^{(i)} + \phi \sum_{i=1}^m (1 - 2y^{(i)}) = 0 \Rightarrow \phi = \frac{\sum_{i=1}^m y^{(i)}}{\sum_{i=1}^m (2y^{(i)} - 1)}$$

$$\phi = \frac{\sum_{i=1}^m \mathbb{1}\{y^{(i)} = 1\}}{m}$$

$\sum_{i=1}^m \mathbb{1}$ [for any $y^{(i)}$]

$$\begin{aligned} l &= \frac{\partial}{\partial \mu_0} \sum_{i=1}^m \log \left(\mathbb{1}\{y^{(i)} = 0\} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x^{(i)} - \mu_0)^2}{2\sigma^2}\right) \right) \\ &= \sum_{i=1}^m \mathbb{1}\{y^{(i)} = 0\} \frac{\partial}{\partial \mu_0} \log \left(e^{-\frac{(x^{(i)} - \mu_0)^2}{2\sigma^2}} \right) \\ &= \sum_{i=1}^m \mathbb{1}\{y^{(i)} = 0\} \frac{\partial}{\partial \mu_0} \left(-\frac{(x^{(i)} - \mu_0)^2}{2\sigma^2} \right) = -\mathbb{1}\{y^{(i)} = 0\} \sum_{i=1}^m \frac{x^{(i)} - \mu_0}{\sigma^2} = 0 \end{aligned}$$

$$\frac{d}{dx} \log(e^x) = e^x \frac{d}{dx} x$$

$\log(e^x) = x$ [log and exp inverse function]

$$S_1 - \sum_{i=1}^m \mathbb{1}\{y^{(i)} = 0\} x^{(i)} + \sum_{i=1}^m \mathbb{1}\{y^{(i)} = 0\} \mu_0 = 0$$

$$\therefore \mu_0 = \sum_{i=1}^m \mathbb{1}\{y^{(i)} = 0\} x^{(i)} / \sum_{i=1}^m \mathbb{1}\{y^{(i)} = 0\}$$

Finally found

$$\text{Note } y = \sigma^2$$

$$\text{Similarly for } \Sigma = \sigma^2, \quad \sqrt{y} = \sigma$$

Bad approximation of GDA on data set 1 \Rightarrow Since the dataset is not Gaussian
Fix (take $\log(x^{(i)})$)

⑦ Implementation Notes

Algorithm

① Compute $\phi, \mu_0, \mu_1, \Sigma$ using above formula

② Compute θ using \ast

③ Find $P(y=1 | x)$; if $P(y=1 | x) \geq 0.5$
predict that label is '1' else predict
that label is 0.

IncompleteLabels

Wednesday, August 12, 2020 1:46 PM

2. [30 points] Incomplete, Positive-Only Labels

In this problem we will consider training binary classifiers in situations where we do not have full access to the labels. In particular, we consider a scenario, which is not too infrequent in real life, where we have labels only for a subset of the positive examples. All the negative examples and the rest of the positive examples are unlabelled.

That is, we assume a dataset $\{(x^{(i)}, t^{(i)}, y^{(i)})\}_{i=1}^m$, where $t^{(i)} \in \{0, 1\}$ is the “true” label, and where

$$y^{(i)} = \begin{cases} 1 & x^{(i)} \text{ is labeled} \\ 0 & \text{otherwise.} \end{cases}$$

All labeled examples are positive, which is to say $p(t^{(i)} = 1 | y^{(i)} = 1) = 1$, but unlabeled examples may be positive or negative. Our goal in the problem is to construct a binary classifier h of the true label t , with only access to the partial labels y . In other words, we want to construct h such that $h(x^{(i)}) \approx p(t^{(i)} = 1 | x^{(i)})$ as closely as possible, using only x and y .

Real world example: Suppose we maintain a database of proteins which are involved in transmitting signals across membranes. Every example added to the database is involved in a signaling process, but there are many proteins involved in cross-membrane signaling which are missing from the database. It would be useful to train a classifier to identify proteins that should be added to the database. In our notation, each example $x^{(i)}$ corresponds to a protein, $y^{(i)} = 1$ if the protein is in the database and 0 otherwise, and $t^{(i)} = 1$ if the protein is involved in a cross-membrane signaling process and thus should be added to the database, and 0 otherwise.

(a) [5 points] Suppose that each $y^{(i)}$ and $x^{(i)}$ are conditionally independent given $t^{(i)}$:

$$p(y^{(i)} = 1 | t^{(i)} = 1, x^{(i)}) = p(y^{(i)} = 1 | t^{(i)} = 1).$$

Note this is equivalent to saying that labeled examples were selected uniformly at random from the set of positive examples. Prove that the probability of an example being labeled differs by a constant factor from the probability of an example being positive. That is, show that $p(t^{(i)} = 1 | x^{(i)}) = p(y^{(i)} = 1 | x^{(i)})/\alpha$ for some $\alpha \in \mathbb{R}$.

Answer:

Qn002

$$\text{a } P(y=1 | t=1, x) P(t=1 | x) P(x) = P(y=1, t=1, x) = P(t=1 | y=1, x) P(y=1 | x) P(x)$$

$$\text{so, } P(t=1 | x) = P(y=1 | x) \frac{P(t=1 | y=1, x)}{P(y=1 | t=1, x)}$$

$$P(t=1 | y=1, x) = 1, P(y=1 | t=1, x) = P(y=1 | t=1)$$

$$P(t=1 | x) = \frac{P(y=1 | x)}{P(y=1 | t=1)} \quad P(y=1 | t=1) = \alpha$$

$$(b) h(x) \approx P(y=1 | x) = P(t=1 | x) \propto \alpha \text{ for all } x \in V,$$

- (b) [5 points] Suppose we want to estimate α using a trained classifier h and a held-out validation set V . Let V_+ be the set of labeled (and hence positive) examples in V , given by $V_+ = \{x^{(i)} \in V \mid y^{(i)} = 1\}$. Assuming that $h(x^{(i)}) \approx p(y^{(i)} = 1 \mid x^{(i)})$ for all examples $x^{(i)}$, show that

$$h(x^{(i)}) \approx \alpha \quad \text{for all } x^{(i)} \in V_+.$$

You may assume that $p(t^{(i)} = 1 \mid x^{(i)}) \approx 1$ when $x^{(i)} \in V_+$.

Answer:

- (c) [5 points] **Coding problem.** The following three problems will deal with a dataset which we have provided in the following files:

`data/ds3_{train,valid,test}.csv`

Each file contains the following columns: x_1 , x_2 , y , and t . As in Problem 1, there is one example per row.

First we will consider the ideal case, where we have access to the true t -labels for training. In `src/p02cde_posonly`, write a logistic regression classifier that uses x_1 and x_2 as input features, and train it using the t -labels (you can ignore the y -labels for this part). Output the trained model's predictions on the test set to the file specified in the code.

Answer:

- (d) [5 points] **Coding problem.** We now consider the case where the t -labels are unavailable, so you only have access to the y -labels at training time. Add to your code in `p02cde_posonly.py` to re-train the classifier (still using x_1 and x_2 as input features), but using the y -labels only.

Answer:

- (e) [10 points] **Coding problem.** Using the validation set, estimate the constant α by averaging your classifier's predictions over all labeled examples in the validation set:

$$\alpha \approx \frac{1}{|V_+|} \sum_{x^{(i)} \in V_+} h(x^{(i)}).$$

Add code in `src/p02cde_posonly.py` to rescale your classifier's predictions from part (d) using the estimated value for α .

Finally, using a threshold of $p(t^{(i)} = 1 \mid x^{(i)}) = 0.5$, make three separate plots with the decision boundaries from parts (c) - (e) plotted on top of the test set. Plot x_1 on the horizontal axis and x_2 on the vertical axis, and use two different symbols for the positive ($t^{(i)} = 1$) and negative ($t^{(i)} = 0$) examples. In each plot, indicate the separating hyperplane with a red line.

Answer:

Remark: We saw that the true probability $p(t \mid x)$ was only a constant factor away from $p(y \mid x)$. This means, if our task is to only rank examples (*i.e.* sort them) in a particular order (e.g., sort the proteins in order of being most likely to be involved in transmitting signals across membranes), then in fact we do not even need to estimate α . The rank based on $p(y \mid x)$ will agree with the rank based on $p(t \mid x)$.

Possion_GLM

Wednesday, August 12, 2020 2:10 PM

3. [25 points] Poisson Regression

- (a) [5 points] Consider the Poisson distribution parameterized by λ :

$$p(y; \lambda) = \frac{e^{-\lambda} \lambda^y}{y!}.$$

Show that the Poisson distribution is in the exponential family, and clearly state the values for $b(y)$, η , $T(y)$, and $a(\eta)$.

Answer:

$$\begin{aligned} p(y; \lambda) &= b(y) \exp(\eta^T T(y) - a(\eta)) \quad (a) \\ p(y; \lambda) &= \frac{e^{-\lambda} \lambda^y}{y!} = \frac{1}{y!} e^{-\lambda} e^{y \log \lambda} = \frac{1}{y!} e^{y \log \lambda - \lambda} \\ \text{So, } b(y) &= 1/y! \quad T(y) = y \\ \eta &= \log \lambda \quad a(\eta) = \lambda = e^\eta \end{aligned}$$

- (b) [3 points] Consider performing regression using a GLM model with a Poisson response variable. What is the canonical response function for the family? (You may use the fact that a Poisson random variable with parameter λ has mean λ .)

Answer:

- (c) [7 points] For a training set $\{(x^{(i)}, y^{(i)}); i = 1, \dots, m\}$, let the log-likelihood of an example be $\log p(y^{(i)}|x^{(i)}; \theta)$. By taking the derivative of the log-likelihood with respect to θ_j , derive the stochastic gradient ascent update rule for learning using a GLM model with Poisson responses y and the canonical response function.

Answer:

$$\begin{aligned} h(\lambda) &= E(T(y)|\lambda) \stackrel{\text{def}}{=} \lambda = e^{\theta^T x} \quad \text{So, } h_\theta(\lambda) = e^{\theta^T x} \left[\begin{array}{c} \delta(\lambda) = e^\lambda \\ \text{canon respn} \end{array} \right] \\ l(\theta) &= \sum_{i=1}^m h(\lambda) e^{\lambda^T x^{(i)} - a(\lambda)} = \sum_{i=1}^m \log \left[\frac{1}{y!} e^{y \log \lambda - \lambda} \right] = \sum_{i=1}^m \log \left(\frac{1}{y!} \right) + \log e^{y \log \lambda} \\ &= \sum_{i=1}^m \log \left(\frac{1}{y!} \right) + y \log \lambda - e^{\theta^T x} \\ \frac{\partial l}{\partial \theta_j} &= \sum_{i=1}^m \frac{\partial}{\partial \theta_j} (\log \left(\frac{1}{y!} \right) + y \log \lambda - e^{\theta^T x}) = \sum_{i=1}^m y^{(i)} x_j^{(i)} - e^{\theta^T x} x_j^{(i)} \\ \text{So, Stochastic Gradient descent rule is: } &\theta := \theta + \alpha \left(y^{(i)} - e^{\theta^T x^{(i)}} \right) x^{(i)} \end{aligned}$$

- (d) [7 points] **Coding problem.** Consider a website that wants to predict its daily traffic. The website owners have collected a dataset of past traffic to their website, along with some features which they think are useful in predicting the number of visitors per day. The dataset is split into train/valid/test sets and follows the same format as Datasets 1-3:

```
data/ds4_{train,valid}.csv
```

We will apply Poisson regression to model the number of visitors per day. Note that applying Poisson regression in particular assumes that the data follows a Poisson distribution whose natural parameter is a linear combination of the input features (*i.e.*, $\eta = \theta^T x$). In `src/p03d.poisson.py`, implement Poisson regression for this dataset and use gradient ascent to maximize the log-likelihood of θ .

Answer:

Convexity of Generalized Linear Models

Wednesday, August 12, 2020 2:21 PM

- (a) [5 points] Derive an expression for the mean of the distribution. Show that $\mathbb{E}[Y | X; \theta]$ can be represented as the gradient of the log-partition function a with respect to the natural parameter η .

Hint: Start with observing that $\frac{\partial}{\partial \eta} \int p(y; \eta) dy = \int \frac{\partial}{\partial \eta} p(y; \eta) dy$.

Answer:

$$\left| \begin{array}{l}
 \int p(y; \eta) dy = 1 \\
 \frac{\partial}{\partial \eta} \int p(y; \eta) dy = 0 \\
 \int \frac{\partial}{\partial \eta} p(y; \eta) dy = 0
 \end{array} \right| \quad \left| \begin{array}{l}
 \int \frac{\partial}{\partial \eta} b(y) \exp(\eta y - a(\eta)) dy = 0 \\
 \int b(y) \exp(\eta y - a(\eta)) (y - a'(\eta)) dy = 0 \\
 \int (g p(y; \eta) - p y \ln a'(\eta)) dy = 0
 \end{array} \right. \\
 \int y p(y; \eta) dy = \int p(y; \eta) a'(\eta) dy \\
 E[Y | x; \eta] = a'(\eta) \int p(y; \eta) dy \\
 \text{so, } E[Y | x; \eta] = a'(\eta)$$

- (b) [5 points] Next, derive an expression for the variance of the distribution. In particular, show that $\text{Var}(Y | X; \theta)$ can be expressed as the derivative of the mean w.r.t η (i.e., the second derivative of the log-partition function $a(\eta)$ w.r.t the natural parameter η .)

Answer:

$$\begin{aligned}
 \frac{\partial^2}{\partial \eta^2} \int p(y; \eta) dy &= \frac{\partial}{\partial \eta} \int (y - a'(\eta)) b(y) \exp(\eta y - a(\eta)) dy \\
 &= \int \frac{\partial}{\partial \eta} (()) dy
 \end{aligned}$$

$$\begin{aligned}
&= \int (y - a'(n))^2 b(y) \exp(ny - an) - b(y) \exp(ny - an) a''(n) \\
&= \int b(y) \exp(ny - a(n)) [(y - a'(n))^2 - a''(n)] dy \\
&= \int p(y; n) [(y - a'(n))^2 - a''(n)] dy \\
&= \int p(y; n) (y - a'(n))^2 dy - \int p(y; n) a''(n) dy = 0
\end{aligned}$$

or $\int p(y; n) [y^2 - 2y a'(n) + a'(n)^2] dy = a''(n)$

or $\int p(y; n) y^2 dy - 2 \underbrace{E[Y|x]}_{E[Y|x]^2} + E[Y|x]^2$

$$E[Y^2|x] - E[Y|x]^2 = a''(n)$$

$$\therefore \text{Var}(\varphi|x; \theta) = a''(n)$$

(c) [5 points] Finally, write out the loss function $\ell(\theta)$, the NLL of the distribution, as a function of θ . Then, calculate the Hessian of the loss w.r.t θ , and show that it is always PSD. This concludes the proof that NLL loss of GLM is convex.

Hint: Use the chain rule of calculus along with the results of the previous parts to simplify your derivations.

Answer:

Remark: The main takeaways from this problem are:

- Any GLM model is convex in its model parameters.
- The exponential family of probability distributions are mathematically nice. Whereas calculating mean and variance of distributions in general involves integrals (hard), surprisingly we can calculate them using derivatives (easy) for exponential family.

$$\begin{aligned}
\ell(\theta) &= \sum_{i=1}^m \log (p(y^{(i)} | x^{(i)}; \theta)) = \sum_{i=1}^m \log p(y^{(i)}, n) \\
\ell(\theta) &= -\sum_{i=1}^m \log b(y^{(i)}) + y^{(i)} \theta^\top x^{(i)} - q(\theta^\top x^{(i)}) \\
\frac{\partial \ell}{\partial \theta_j} &= \frac{\partial}{\partial \theta_j} \ell(\theta) = -\frac{\partial}{\partial \theta_j} \left(\sum_{i=1}^m y^{(i)} x_j^{(i)} - a'(\theta^\top x^{(i)}) x_k^{(i)} \right)
\end{aligned}$$

$$= \sum_{i=1}^m \left[a''(\theta^T x^{(i)}) x_j^{(i)} x_k^{(i)} \right]$$

∴

$$\begin{aligned} z^T H z &= \sum_i \sum_j z_i h_{ij} z_j = \sum_i \sum_j z_i a''(\theta^T x^{(i)}) x_j x_k z_j \\ &= a''(\theta^T x) (\theta^T z)^2 \end{aligned}$$

Since, $a''(\theta^T x)$ is $\text{Var}(Y|x; \theta)$ so, $a''(\theta^T x)$ is non-negative

Hence $z^T H z \geq 0 \Rightarrow l$ is convex

Locally weighted linear regression

Wednesday, August 12, 2020 2:29 PM

5. [25 points] Locally weighted linear regression

- (a) [10 points] Consider a linear regression problem in which we want to "weight" different training examples differently. Specifically, suppose we want to minimize

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m w^{(i)} (\theta^T x^{(i)} - y^{(i)})^2.$$

In class, we worked out what happens for the case where all the weights (the $w^{(i)}$'s) are the same. In this problem, we will generalize some of those ideas to the weighted setting.

- i. [2 points] Show that $J(\theta)$ can also be written

$$J(\theta) = (X\theta - y)^T W (X\theta - y)$$

for an appropriate matrix W , and where X and y are as defined in class. Clearly specify the value of each element of the matrix W .

- ii. [4 points] If all the $w^{(i)}$'s equal 1, then we saw in class that the normal equation is

$$X^T X \theta = X^T y,$$

and that the value of θ that minimizes $J(\theta)$ is given by $(X^T X)^{-1} X^T y$. By finding the derivative $\nabla_{\theta} J(\theta)$ and setting that to zero, generalize the normal equation to this weighted setting, and give the new value of θ that minimizes $J(\theta)$ in closed form as a function of X , W and y .

- iii. [4 points] Suppose we have a dataset $\{(x^{(i)}, y^{(i)}); i = 1 \dots, m\}$ of m independent examples, but we model the $y^{(i)}$'s as drawn from conditional distributions with different levels of variance $(\sigma^{(i)})^2$. Specifically, assume the model

$$p(y^{(i)} | x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi\sigma^{(i)}}} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2(\sigma^{(i)})^2}\right)$$

That is, each $y^{(i)}$ is drawn from a Gaussian distribution with mean $\theta^T x^{(i)}$ and variance $(\sigma^{(i)})^2$ (where the $\sigma^{(i)}$'s are fixed, known, constants). Show that finding the maximum likelihood estimate of θ reduces to solving a weighted linear regression problem. State clearly what the $w^{(i)}$'s are in terms of the $\sigma^{(i)}$'s.

Answer:

$$(i) X \in \begin{bmatrix} 1 & x_1 & \cdots & x_n \end{bmatrix}^T$$

$$W \in \begin{bmatrix} (1) & & & \\ \frac{w^{(1)}}{2} & \frac{w^{(2)}}{2} & \cdots & 0 \\ 0 & \ddots & \ddots & \frac{w^{(m)}}{2} \end{bmatrix}$$

$$\text{st } / \quad w(X\theta - y) = \begin{bmatrix} \frac{w^{(1)}}{2} & (\theta^T x^{(1)} - y^{(1)}) \\ \vdots & \vdots \\ \frac{w^{(m)}}{2} & (\theta^T x^{(m)} - y^{(m)}) \end{bmatrix}$$

$$\begin{aligned}
 \nabla_{\theta} J(\theta) &= D_{\theta} (X\theta - y)^T w (X\theta - y) \\
 &= D_{\theta} (\theta^T x^T - y^T) (w x \theta - w y) \\
 &= D_{\theta} (\theta^T x^T w x \theta - \theta^T x^T w y - y^T w x \theta + y^T w y) \\
 &= D_{\theta} (\theta^T x^T w x \theta - 2 y^T w x \theta + y^T w y) \\
 &= 2 x^T w x \theta - 2 x^T w^T y = 0
 \end{aligned}$$

$$x^T w x \theta = x^T w^\dagger y$$

$$\theta = (x^T w x)^{-1} x^T w y$$

ii) $\ell(\theta) = \sum_{i=1}^m -\log \tau_2 \sigma^{(i)} - \log \sigma^{(i)} - \frac{1}{2\sigma^{(i)2}} (y^{(i)} - \theta^T x^{(i)})^2$

Maximizing $\ell(\theta)$ is equivalent to minimizing $J\theta = \sum_{i=1}^m \frac{1}{2\sigma^{(i)2}} (y^{(i)} - \theta^T x^{(i)})^2$

$$\text{so, } w^{(i)} = \frac{1}{(\sigma^{(i)})^2}$$

- (b) [10 points] **Coding problem.** We will now consider the following dataset (the formatting matches that of Datasets 1-4, except $x^{(i)}$ is 1-dimensional):

`data/ds5_{train,valid,test}.csv`

In `src/p05b_lwr.py`, implement locally weighted linear regression using the normal equations you derived in Part (a) and using

$$w^{(i)} = \exp\left(-\frac{\|x^{(i)} - x\|_2^2}{2\tau^2}\right).$$

Train your model on the `train` split using $\tau = 0.5$, then run your model on the `valid` split and report the mean squared error (MSE). Finally plot your model's predictions on the validation set (plot the training set with blue 'x' markers and the validation set with a red 'o' marker). Does the model seem to be under- or overfitting?

Answer:

- (c) [5 points] **Coding problem.** We will now tune the hyperparameter τ . In `src/p05c_tau.py`, find the MSE value of your model on the validation set for each of the values of τ specified in the code. For each τ , plot your model's predictions on the validation set in the format described in part (b). Report the value of τ which achieves the lowest MSE on the `valid` split, and finally report the MSE on the `test` split using this τ -value.

Answer: