

$$\log e = x$$

Problem 1: MLE of Multivariate Gaussian Distribution

$$x_1, \dots, x_n \in \mathbb{R}^d$$

$$x \sim N(\mu, \sigma^2 I_{\text{nd}})$$

$$P(x_i) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} e^{-\frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)}$$

$$L = \sum_{i=1}^n \log P(x_i)$$

$$\log P(x_i) = -\frac{d}{2} \log 2\pi - \frac{1}{2} \log(|\Sigma|) - \frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)$$

$$L = -\frac{Nd}{2} \log 2\pi - \frac{N}{2} \log |\Sigma| - \sum_{i=1}^N \frac{(x_i - \mu)^T \Sigma^{-1} (x_i - \mu)}{2}$$

$$\Sigma = \begin{bmatrix} \sigma^2 & & \\ & \sigma^2 & \\ & & \ddots \\ & & & \sigma^2 \end{bmatrix} \rightarrow |\Sigma| = (\sigma^2)(\sigma^2) \dots (\sigma^2) = \sigma^{2d}$$

$$\Sigma^{-1} = \begin{bmatrix} \sigma^{-2} & & \\ & \sigma^{-2} & \\ & & \ddots \\ & & & \sigma^{-2} \end{bmatrix} \rightarrow (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) = \sigma^{-2} (x_i - \mu)^T (x_i - \mu)$$

$$L = -\frac{nd}{2} \log(2\pi) - nd \log \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^T (x_i - \mu)$$

$$L = nd(-\log \sqrt{2\pi} - \log \sigma) - \sum_{i=1}^n \frac{\|x_i - \mu\|_2^2}{2\sigma^2}$$

$$\frac{\partial L}{\partial \mu} = -\frac{1}{2\sigma^2} (-2) \sum_{i=1}^n (x_i - \mu) = 0 \Rightarrow \sum_{i=1}^n x_i - n\mu = 0$$

$$\mu = \frac{\sum_{i=1}^n x_i}{n}$$

$$\frac{\partial L}{\partial \sigma} = -\frac{nd}{\sigma} - \frac{1}{\sigma^3} \sum_{i=1}^n (x_i - \mu)^T (x_i - \mu) = 0$$

$$0 = -\frac{nd}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^n \|x_i - \mu\|_2^2 \Rightarrow \sigma^2 = \frac{\sum_{i=1}^n \|x_i - \mu\|_2^2}{nd}$$

$$\{1, 2, 3, 4, 5\} \{10\} \begin{bmatrix} 5 \\ 4 \\ 5 \end{bmatrix} 10$$

2 b) $X \sim N(\mu, \Lambda)$ where $\Lambda = \begin{bmatrix} \sigma_1^2 & \sigma_2^2 & \dots & \sigma_d^2 \end{bmatrix}$ $X: x_1, \dots, x_N \in \mathbb{R}^d$

$$P(x_i) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} e^{-\frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)}$$

$$L = \sum_{i=1}^N \log P(x_i)$$

$$\log P(x_i) = -d \log \sqrt{2\pi} - \log \sqrt{|\Sigma|} - \frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)$$

$$|\Sigma| = \sigma_1^2 \sigma_2^2 \dots \sigma_d^2 = \left(\prod_{j=1}^d \sigma_j \right)^2 \Rightarrow \sqrt{|\Sigma|} = \prod_{j=1}^d \sigma_j$$

$$\log \sqrt{|\Sigma|} = - \sum_{j=1}^d \log \sigma_j$$

Trace method \rightarrow

$$\Sigma^{-1} = \begin{bmatrix} \sigma_1^{-2} & & \\ & \sigma_2^{-2} & \\ & & \ddots \\ & & & \sigma_d^{-2} \end{bmatrix} \Rightarrow \frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) = \frac{1}{2} [\sigma_1^{-2} \sigma_2^{-2} \dots \sigma_d^{-2}] \|x_i - \mu\|_2^2$$

$$= \frac{1}{2} \text{Tr}(\Sigma^{-1} \|x_i - \mu\|_2^2)$$

$$= \frac{1}{2} \|x_i - \mu\|_2^2 \sum_{j=1}^d \sigma_j^{-2}$$

$$L = \log P(x_i) = -d \log \sqrt{2\pi} - \sum_{j=1}^d \log \sigma_j - \frac{1}{2} \|x_i - \mu\|_2^2 \sum_{j=1}^d \sigma_j^{-2}$$

$$L = -Nd \log \sqrt{2\pi} - N \sum_{j=1}^d \log \sigma_j - \sum_{j=1}^d \frac{1}{2\sigma_j^2} \sum_{i=1}^N \|x_i - \mu\|_2^2$$

$$\frac{\partial L}{\partial \mu} = - \sum_{j=1}^d \frac{1}{2\sigma_j^2} \sum_{i=1}^N (x_i - \mu) = 0 = (\sigma_1^{-2} + \sigma_2^{-2} + \dots + \sigma_d^{-2}) \sum_{i=1}^N (x_i - \mu) = 0$$

$$\mu = \frac{\sum_{i=1}^N x_i}{N}$$

$$\frac{\partial L}{\partial \Lambda} = -N \sum_{j=1}^d \frac{1}{\sigma_j} + \sum_{j=1}^d \frac{1}{2\sigma_j^3} \sum_{i=1}^N \|x_i - \mu\|_2^2 = 0$$

$$N (\sigma_1^{-1} + \sigma_2^{-1} + \dots + \sigma_d^{-1}) = (\sigma_1^{-3} + \sigma_2^{-3} + \dots + \sigma_d^{-3}) \sum_{i=1}^N \|x_i - \mu\|_2^2$$

$$\sum_{j=1}^d \frac{1}{\sigma_j} = \frac{\sum_{i=1}^N \|x_i - \mu\|_2^2}{\sum_{j=1}^d \left(\frac{1}{\sigma_j}\right)^3} = \frac{\sum_{i=1}^N \|x_i - \mu\|_2^2}{N}$$

$$X \sim N(\mu, \Lambda) \text{ where } \Lambda = \begin{bmatrix} \sigma_1^2 & & \\ & \sigma_2^2 & \\ & & \ddots \\ & & & \sigma_d^2 \end{bmatrix}$$

$$\log P(x_i) = -d \log \sqrt{2\pi} - \log \sqrt{|\Lambda|} - \frac{1}{2} (x_i - \mu)^T \Lambda^{-1} (x_i - \mu)$$

$$|\Lambda| = \text{Tr } \Lambda \Rightarrow -\log \sqrt{|\Lambda|} = -\frac{1}{2} \log \text{Tr } \Lambda$$

$$L = \sum_{i=1}^N \log P(x_i) = -Nd \log \sqrt{2\pi} - \frac{N}{2} \log (\text{Tr } \Lambda) - \sum_{i=1}^N \frac{(x_i - \mu)^T \Lambda^{-1} (x_i - \mu)}{2}$$

$$\frac{\partial L}{\partial \mu} = + \sum_{i=1}^N \frac{\Lambda^{-1} (x_i - \mu)}{2} = 0 \Rightarrow \boxed{\mu = \frac{\sum_{i=1}^N x_i}{N}}$$

$$\frac{\partial L}{\partial \Lambda} = \frac{-N}{2} \frac{1}{\text{Tr } \Lambda} \frac{\partial \text{Tr } \Lambda}{\partial \Lambda} - \frac{1}{2} \sum_{i=1}^N \frac{\partial}{\partial \Lambda} (x_i - \mu)^T \Lambda^{-1} (x_i - \mu) = 0$$

$$L = -Nd \log \sqrt{2\pi} - \frac{N}{2} \log |\Lambda| - \frac{1}{2} \sum_{i=1}^N (x_i - \mu)^T \Lambda^{-1} (x_i - \mu)$$

$$\frac{\partial L}{\partial \Lambda} = \frac{-N}{2} \frac{1}{|\Lambda|} \frac{\partial |\Lambda|}{\partial \Lambda} - \frac{1}{2} \sum_{i=1}^N \frac{\partial}{\partial \Lambda} (x_i - \mu)^T \Lambda^{-1} (x_i - \mu)$$

$$\frac{\partial |\Lambda|}{\partial \Lambda} = |\Lambda| (\Lambda^{-1})^T = |\Lambda| (\Lambda^T)^{-1} = |\Lambda| \Lambda^{-1}$$

$$\frac{\partial L}{\partial \Lambda} = \frac{-N}{2} (\Lambda^{-1})^T - \frac{1}{2} \sum_{i=1}^N \frac{\partial}{\partial \Lambda} (x_i - \mu)^T \Lambda^{-1} (x_i - \mu) = 0$$

$$\Lambda \Lambda^{-1} = I \Rightarrow \frac{\partial}{\partial \Lambda} (\Lambda \Lambda^{-1}) = 0 = \left(\frac{\partial}{\partial \Lambda} \Lambda \right) \Lambda^{-1} + \Lambda \left(\frac{\partial}{\partial \Lambda} \Lambda^{-1} \right) = 0$$

$$\frac{\partial}{\partial \Lambda} \Lambda^{-1} = -\Lambda^{-1} \left(\frac{\partial}{\partial \Lambda} \Lambda \right) \Lambda^{-1} = -\Lambda^{-1} \Lambda^{-1}$$

$$\text{b/c } \Lambda^{-1} \text{ diagonal} \Rightarrow \frac{\partial}{\partial \Lambda} \Lambda^{-1} = -\Lambda^{-2}$$

$$\frac{\partial L}{\partial \Lambda} = 0 = \frac{-N \Lambda^{-1}}{2} + \frac{1}{2} \sum_{i=1}^N \|x_i - \mu\|_2^2 \Lambda^{-2} = 0$$

$$\boxed{\Lambda = \frac{\sum_{i=1}^N \|x_i - \mu\|_2^2}{N}}$$

c). $X \sim N(A\mu, \Lambda)$ $A, D \in \mathbb{R}^{d \times d}$ $\mu = ?$

$$P(x_i) = \frac{1}{\sqrt{(2\pi)^d |\Lambda|}} e^{-\frac{1}{2} (x_i - A\mu)^T \Lambda^{-1} (x_i - A\mu)}$$

$$\log P(x_i) = -d \log \sqrt{2\pi} - \log \sqrt{|\Lambda|} - \frac{1}{2} (x_i - A\mu)^T \Lambda^{-1} (x_i - A\mu)$$

$$L = \sum \log P(x_i) = -Nd \log \sqrt{2\pi} - N \log \sqrt{|\Lambda|} - \frac{N}{2} \sum (x_i - A\mu)^T \Lambda^{-1} (x_i - A\mu)$$

$$\frac{\partial L}{\partial \mu} = \frac{N}{2} \sum_{i=1}^N \frac{\partial}{\partial \mu} (x_i - A\mu)^T \Lambda^{-1} (x_i - A\mu) = 0$$

$$= \frac{N}{2} \sum_{i=1}^N -2A \Lambda^{-1} (x_i - A\mu) = 0$$

$$= -A \Lambda^{-1} \sum_{i=1}^N (x_i - A\mu) = 0$$

$$\boxed{\mu = \frac{\sum_{i=1}^N x_i}{nA}}$$

$$(x+\mu)^T(x+\mu) = \|x-\mu\|^2$$

Problem 2: ℓ_2 -regularized logistic/linear Regression w/ Newton's Method

$S = \{(x_i, y_i)\}_{i=1}^n$ training set $x_i \in \mathbb{R}, y_i \in \{0, 1\}$

a). Negative log-likelihood for ℓ_2 -regularized logistic regression:

$$l(\beta) = \lambda \|\beta\|_2^2 - \sum_{i=1}^n [y_i \log \mu_i + (1-y_i) \log(1-\mu_i)]$$

$$\mu_i = \frac{1}{1 + e^{-\beta x_i}} \quad \lambda > 0 \text{ is regularization parameter}$$

b). ℓ_2 -regularized quadratic cost function (i.e. ridge regression)

$$J(\beta) = \lambda \|\beta\|_2^2 + \frac{1}{2} \sum_{i=1}^n (y_i - \beta x_i)^2$$

$$X = \begin{bmatrix} 4 & 1 \\ 5 & 1 \\ 5.0 & 1 \\ \vdots & \vdots \\ 4.3 & 1 \end{bmatrix} \begin{matrix} \leftarrow x_1 \\ \leftarrow x_2 \\ \\ \leftarrow x_n \end{matrix} \quad B = \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ 1 \\ 1 \\ \vdots \\ 0 \end{bmatrix} \begin{matrix} \leftarrow y_1 \\ \leftarrow y_2 \\ \\ \leftarrow y_n \end{matrix} \quad M = \begin{bmatrix} \frac{1}{1+e^{-\beta x_1}} \\ \frac{1}{1+e^{-\beta x_2}} \\ \vdots \\ \frac{1}{1+e^{-\beta x_n}} \end{bmatrix}$$

1. $\nabla_{\beta} l(\beta) = ?$, $\nabla_{\beta} J(\beta) = ?$

$$l(\beta) = \lambda \|\beta\|_2^2 - Y^T \log M + (1 - Y^T) \log(1 - M)$$

$$\nabla l(\beta) = 2\lambda\beta - Y^T \log(1 + e^{-\beta X})^{-1} + (1 - Y^T) \log(1 - (1 + e^{-\beta X})^{-1})$$

$$\begin{aligned} l(\beta) &= \lambda \|\beta\|_2^2 + Y^T \log(1 + e^{-\beta X}) + (1 - Y^T) \log((2 + e^{-\beta X})(1 + e^{-\beta X})^{-1}) \\ &= \lambda \|\beta\|_2^2 + Y^T \log(1 + e^{-\beta X}) + (1 - Y^T) \log(2 + e^{-\beta X}) - (1 - Y^T) \log(1 + e^{-\beta X}) \\ &= \lambda \|\beta\|_2^2 + (2Y^T - 1) \log(1 + e^{-\beta X}) + (1 - Y^T) \log(2 + e^{-\beta X}) \end{aligned}$$

$$\nabla_{\beta} l(\beta) = 2\lambda\beta + \frac{2Y^T - 1}{1 + e^{-\beta X}} (-X e^{-\beta X}) + \frac{1 - Y^T}{2 + e^{-\beta X}} (-X e^{-\beta X}) =$$

1. $\nabla_{\beta} l(\beta) = ?$, $\nabla_{\beta} J(\beta) = ?$

$$\nabla_{\beta} l(\beta) = 2\lambda\beta - \sum_{i=1}^n \frac{y_i}{\mu_i} \frac{\partial \mu_i}{\partial \beta} + \frac{1-y_i}{1-\mu_i} \frac{\partial \mu_i}{\partial \beta} = 0$$

$$\frac{\partial \mu_i}{\partial \beta} = \frac{\partial}{\partial \beta} \left(\frac{1}{1 + e^{-\beta x_i}} \right) = \frac{1}{(1 + e^{-\beta x_i})^2} (x_i e^{-\beta x_i}) = \frac{x_i e^{-\beta x_i}}{(1 + e^{-\beta x_i})^2}$$

$$\nabla_{\beta} l(\beta) = 0 = 2\lambda\beta - \sum_{i=1}^n \frac{y_i}{\mu_i} \mu_i^2 x_i e^{-\beta x_i} - \frac{1-y_i}{1-\mu_i} \mu_i^2 x_i e^{-\beta x_i}$$

$$0 = 2\lambda\beta - \sum_{i=1}^n \frac{(2y_i - 1) x_i e^{-\beta x_i}}{1 + e^{-\beta x_i}} \Rightarrow 2\lambda\beta - \frac{(2Y^T - 1) X e^{-\beta X}}{(1 - X\beta)} = 0$$

1. Find $\nabla_{\beta} \ell(\beta)$ $\nabla_{\beta} J(\beta)$

a) $\nabla_{\beta} \ell(\beta) = 2\lambda\beta - \sum_{i=1}^n \frac{y_i}{\mu_i} \frac{\partial \mu_i}{\partial \beta} - \frac{1-y_i}{1-\mu_i} \frac{\partial \mu_i}{\partial \beta} = 0$

$$\frac{\partial \mu_i}{\partial \beta} = \frac{-1}{(1+e^{-x_i\beta})^2} \cdot (-x_i e^{-x_i\beta}) = \frac{-1}{1+e^{-x_i\beta}} \cdot \frac{e^{-x_i\beta}}{1+e^{-x_i\beta}} \cdot x_i$$

$$= \mu_i (1-\mu_i) x_i$$

$\nabla_{\beta} \ell(\beta) = 0 = 2\lambda\beta - \sum_{i=1}^n \left(\frac{y_i}{\mu_i} - \frac{1-y_i}{1-\mu_i} \right) \mu_i (1-\mu_i) x_i$

$$2\lambda\beta = \sum_{i=1}^n (y_i(1-\mu_i) - (1-y_i)\mu_i) x_i$$

$$= \sum_{i=1}^n (y_i - \mu_i y_i - \mu_i + \mu_i y_i) x_i$$

$$2\lambda\beta = \sum_{i=1}^n (y_i - \mu_i) x_i$$

$$\beta = \frac{\sum_{i=1}^n (y_i - \mu_i) x_i}{2\lambda} = \frac{1}{2\lambda} \frac{X^T(Y-M)}{2\lambda}$$

$$\boxed{\nabla_{\beta} \ell(\beta) = 2\lambda\beta - X^T(Y-M)} \quad \boxed{2\lambda}$$

b) $J(\beta) = \lambda \|\beta\|_2^2 + \frac{1}{2} (Y - X\beta)^T (Y - X\beta)$

$$\nabla J(\beta) = 2\lambda\beta + \frac{1}{2} (Y^T - \beta^T X^T) (Y - X\beta)$$

$$J = \lambda \|\beta\|_2^2 + \frac{1}{2} (Y^T Y - 2Y^T X\beta + \beta^T X^T X\beta)$$

$$\nabla_{\beta} J(\beta) = 0 = 2\lambda\beta + \frac{1}{2} (-2Y^T X + 2X^T X\beta)$$

$$Y^T X = 2\lambda\beta + X^T X\beta$$

$$\boxed{\beta = (X^T X + 2\lambda I)^{-1} Y^T X}$$

$$\boxed{\nabla_{\beta} J(\beta) = (X^T X + 2\lambda I)\beta - Y^T X}$$

2. a) $\nabla_{\beta}^2 \ell(\beta) = \frac{\partial}{\partial \beta} \nabla_{\beta} \ell(\beta) = \frac{\partial}{\partial \beta} (2\lambda\beta - X^T(Y-M)) = \boxed{2\lambda + X^T M(1-\mu)X}$

b) $\nabla_{\beta}^2 J(\beta) = \frac{\partial}{\partial \beta} \nabla_{\beta} J(\beta) = \frac{\partial}{\partial \beta} (2\lambda\beta + X^T X\beta - Y^T X)$

$$= \boxed{2\lambda I + X^T X}$$

$$3. \quad \ell(\beta) \rightarrow \beta^* \leftarrow \beta - \frac{\nabla_{\beta} \ell(\beta)}{\nabla_{\beta}^2 \ell(\beta)} = \beta - \frac{2\lambda\beta - X^T(Y - \mu)}{2\lambda}$$

$$\boxed{\beta_{k+1} \leftarrow \beta_k - \frac{X^T(Y - \mu)}{2\lambda}}$$

$$J(\beta) \rightarrow \beta_J \leftarrow \beta_J + \frac{\nabla_{\beta} J(\beta)}{\nabla_{\beta}^2 J(\beta)} = \beta_J - \frac{2\lambda\beta + X^T X \beta - Y^T X}{2\lambda I + X^T X}$$

$$\boxed{\beta_{J+1} \leftarrow \beta_J - \frac{Y^T X}{(X^T X + 2\lambda I)^2}}$$

$x_1 \times x_2$

$\mu_1(1-\mu_1)$
 $\mu_2(1-\mu_2)$

$$2. \quad \nabla_{\beta}^2 \ell(\beta) = \frac{\partial}{\partial \beta} \nabla_{\beta} \ell(\beta) = \frac{\partial}{\partial \beta} (2\lambda\beta - X^T(Y - \mu)) = 2\lambda + \frac{\partial}{\partial \beta} (X^T \mu) \\ = 2\lambda + X^T X^T \mu^T (1 - \mu)$$

$$\nabla_{\beta}^2 J(\beta) = \frac{\partial}{\partial \beta} (X^T X + 2\lambda I) \beta - Y^T X = X^T X + 2\lambda I$$

$$3. \quad \ell(\beta) \rightarrow \beta^* \leftarrow \beta - \frac{\nabla_{\beta} \ell(\beta)}{\nabla_{\beta}^2 \ell(\beta)} = \frac{\beta - (2\lambda\beta - X^T(Y - \mu))}{2\lambda I + (X^T)^2 \mu^T (1 - \mu)}$$

$$J(\beta) \rightarrow \beta^* \leftarrow \beta - \frac{\nabla_{\beta} J(\beta)}{\nabla_{\beta}^2 J(\beta)} = \frac{\beta - (X^T X + 2\lambda I) \beta - Y^T X}{2\lambda I + X^T X}$$

$$[\beta_1, \beta_2, \dots, \beta_d]$$

$$\begin{bmatrix} x_1^T y \\ x_2^T y \\ \vdots \\ x_d^T y \end{bmatrix}$$

$$X^T = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_d^T \end{bmatrix}$$

$$X = \begin{bmatrix} | & | & & | \\ x_1 & x_2 & \dots & x_d \\ | & | & & | \end{bmatrix}$$

Problem 3: ℓ_1 -Regularized Linear Regression

$$X \triangleq \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$y \triangleq \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$$\beta \triangleq \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_d \end{bmatrix}$$

Lasso reg

$$\text{want } y = X\beta \quad \beta = \beta^*$$

$$X^T X = nI$$

Lasso Regression

$$\beta^* = \underset{\beta}{\operatorname{argmin}} \{ J_\lambda(\beta) \} = \underset{\beta}{\operatorname{argmin}} \left\{ \frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1 \right\}$$

1. Show β^* determined by i^{th} feature outputs.

$$\text{Write } J_\lambda(\beta) = g(\vec{y}) + \sum_{i=1}^d f(x_i, \vec{y}, \beta_i, \lambda) \quad \text{where } x_i \triangleq i^{\text{th}} \text{ column of } X$$

$$J_\lambda(\beta) = \frac{1}{2} (y - X\beta)^T (y - X\beta) + \lambda \sum_{i=1}^d |\beta_i|$$

$$= \frac{1}{2} (y^T - \beta^T X^T) (y - X\beta) + \lambda \sum_{i=1}^d |\beta_i|$$

$$= \frac{1}{2} (y^T y - 2\beta^T X^T y + \beta^T X^T X \beta) + \lambda \sum_{i=1}^d |\beta_i|$$

$$= \frac{1}{2} y^T y - \sum_{i=1}^d \beta_i X_i^T y + \frac{1}{2} \beta^T nI \beta + \lambda \sum_{i=1}^d |\beta_i|$$

$$= \frac{1}{2} y^T y + \frac{n}{2} \beta^T \beta + \lambda \sum_{i=1}^d |\beta_i| - \beta_i X_i^T y$$

$$J_\lambda(\beta) = \frac{1}{2} (y^T y) + \frac{\lambda n}{2} \sum_{i=1}^d \beta_i^2 + |\beta_i| - \beta_i X_i^T y$$

$$g(\vec{y}) = \frac{1}{2} y^T y$$

$$f(x_i, \vec{y}, \beta_i, \lambda) = \frac{\lambda n}{2} (\beta_i^2 + |\beta_i| - \beta_i X_i^T y)$$

2. Assume $\beta_i^* > 0$ ~~or~~ $J_\lambda(\beta) > 0$

$$\frac{1}{2} (y^T y) + \frac{\lambda n}{2} \sum_{i=1}^d \beta_i^{*2} + |\beta_i^*| - \beta_i^* X_i^T y > 0$$

$$J_\lambda(\beta) = \frac{1}{2} (y^T y) + \frac{\lambda n}{2} \sum_{i=1}^d \beta_i (\beta_i + 1 - X_i^T y)$$

$$\frac{\partial J(\beta)}{\partial \beta_i} = \frac{\lambda n}{2} (2\beta_i + 1 - X_i^T y) = 0$$

$$\beta_i^* = \frac{X_i^T y - 1}{2}$$

3. Assume $\beta_i^* < 0$

$$J(\beta) = \frac{1}{2} y^T y + \frac{\lambda n}{2} \sum_{i=1}^n \beta_i^2 - \beta_i - \beta_i x_i^T y$$

$$\frac{\partial J(\beta)}{\partial \beta_i} = \frac{\lambda n}{2} (2\beta_i - 1 - x_i^T y) = 0$$

$$\boxed{\beta_i^* = \frac{x_i^T y + 1}{2}}$$

(?) (4) When $\alpha = \beta_i^* = 0$

$$\frac{x_i^T y + 1}{2} = 0 \Rightarrow x_i^T y = -1$$

$$\Rightarrow |x_i^T y| = 1$$

$$\frac{x_i^T y - 1}{2} = 0 \Rightarrow x_i^T y = 1$$

5. Lasso: $J_L(\beta) = \frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1$

Ridge: $J_r(\beta) = \frac{1}{2} \|y - X\beta\|_2^2 + \lambda \|\beta\|_2^2$

$$\begin{aligned} J_n(\beta) &= \frac{1}{2} y^T y + \sum_{i=1}^n \frac{\lambda n}{2} (\beta_i^2 + \beta_i^2 - \beta_i x_i^T y) \\ &= \frac{1}{2} y^T y + \sum_{i=1}^n \frac{\lambda n}{2} (2\beta_i^2 - \beta_i x_i^T y) \end{aligned}$$

$$\frac{\partial J_n(\beta)}{\partial \beta_i} = 0 = \frac{\lambda n}{2} (4\beta_i - x_i^T y)$$

$$\boxed{\beta_i^* = \frac{x_i^T y}{4}}$$

$$\Rightarrow \boxed{\beta_i^* = 0 \text{ when } x_i^T y = 0}$$

$$Y \log M + (1-Y) \log(1-M)$$

Problem 4. Span classification using logistic Regression

1. Derive gradient descent equations for logistic regression w/ L2 regularization

$$P(Y=1|X) = \mu_i = \frac{1}{1 + e^{-\beta x_i}} \quad \ell(\beta) = \lambda \|\beta\|_2^2 - \sum_{i=1}^n y_i \log \mu_i + (1-y_i) \log(1-\mu_i)$$

$$\text{where } \mu_i = \frac{1}{1 + e^{-\beta x_i}} \quad \lambda > 0$$

$$\nabla_{\beta} \ell(\beta) = 2\lambda\beta - \sum_{i=1}^n \frac{y_i}{\mu_i} \frac{\partial \mu_i}{\partial \beta} + \frac{1-y_i}{1-\mu_i} \frac{\partial \mu_i}{\partial \beta}$$

$$\frac{\partial \mu_i}{\partial \beta} = \frac{1}{(1 + e^{-\beta x_i})^2} x_i e^{-\beta x_i} = \frac{1}{(1 + e^{-\beta x_i})} \frac{e^{-\beta x_i}}{1 + e^{-\beta x_i}} x_i$$

$$\frac{\partial \mu_i}{\partial \beta} = \mu_i (1 - \mu_i) x_i$$

$$\nabla_{\beta} \ell(\beta) = 2\lambda\beta - \sum_{i=1}^n \left(\frac{y_i}{\mu_i} - \frac{1-y_i}{1-\mu_i} \right) \mu_i (1 - \mu_i) x_i$$

$$= 2\lambda\beta - \sum_{i=1}^n [y_i - y_i \mu_i - (\mu_i - y_i \mu_i)] x_i$$

$$= 2\lambda\beta - \sum_{i=1}^n (y_i - \mu_i) x_i$$

$$\nabla_{\beta} \ell(\beta) = 2\lambda\beta - X^T(y - \mu)$$

$$\beta_{t+1} = \beta_t - \alpha (2\lambda\beta_t - X^T(y - \mu)) = \beta_t - \alpha \sum_{i=1}^n 2\lambda\beta_t - (y_i - \mu_i) x_i$$

See hw3.py

$$2. \quad \beta_{t+1} = \beta_t - \alpha (2\lambda\beta_t - X^T(y - \mu))$$

3. Decreasing alpha over time will allow it to settle into a better optimal solution, so Yes

4. See Haggie

In [1]:

```

import sklearn.metrics as metrics
import numpy as np
import scipy
import scipy.io
import matplotlib.pyplot as plt
import random
import csv
import pdb

%matplotlib inline

clamped = 0.000001
iterations = np.arange(0, 200000, 20000)

def load_dataset():
    mat = scipy.io.loadmat('spam')
    return mat['Xtrain'], mat['Xtest'], mat['ytrain']

def train_gd(X_train, y_train, alpha=1e-7, reg=0.1, num_iter=10000):
    ''' Build a model from X_train -> y_train using batch gradient descent '''
    W = np.zeros((X_train.shape[1], 1))+clamped
    for i in range(num_iter):
        mu = scipy.special.expit(X_train.dot(W))
        gradient = 2*reg*W - X_train.T.dot(y_train - mu)
        W -= alpha * gradient
    return W

def train_sgd(X_train, y_train, alpha=1e-4, reg=0, num_iter=10000):
    ''' Build a model from X_train -> y_train using stochastic gradient descent '''
    W = np.zeros((X_train.shape[1], 1))+clamped
    for i in range(num_iter):
        sample_index = random.randint(0, X_train.shape[0] - 1)
        x_i, y_i = X_train[sample_index].reshape(X_train.shape[1], 1), y_train[sample_index]
        mu = scipy.special.expit(W.T.dot(x_i))
        gradient = 2*reg*W - x_i.dot(y_i - mu)
        W -= alpha * gradient
    return W

def train_sgd_decay(X_train, y_train, alpha=1e-4, reg=0, num_iter=10000):
    ''' Build a model from X_train -> y_train using stochastic gradient descent '''
    W = np.zeros((X_train.shape[1], 1))+clamped
    for i in range(num_iter):
        sample_index = random.randint(0, X_train.shape[0] - 1)
        x_i, y_i = X_train[sample_index].reshape(X_train.shape[1], 1), y_train[sample_index]
        mu = scipy.special.expit(W.T.dot(x_i))
        gradient = 2*reg*W - x_i.dot(y_i - mu)
        W -= (alpha/(i+1)) * gradient
    return W

def predict(model, X):
    ''' From model and data points, output prediction vectors '''

```

```

y_pred = scipy.special.expit(X.dot(model))
return np.round(y_pred)

### Preprocessing Techniques ###
def standardize(X):
    ''' Standardize columns to have mean 0 and unit variance '''
    return (X - np.mean(X, axis=0)) / np.std(X, axis=0)

def log_lift(X):
    ''' Transform the features using log lift '''
    return np.log(X + 0.01)

def binarize(X):
    ''' Binarize features into negative (0) and nonnegative (1) '''
    return np.where(X>0.0, 1.0, 0.0)

def phi(X, G, b):
    ''' Featurize the inputs using random Fourier features '''
    return np.sqrt(2.0/P) * np.cos(np.add(G.dot(X.T), b).T)

def gaussian_random_lift(d):
    mean = np.zeros(d)
    cov = VARIANCE * np.identity(d)
    G = np.random.multivariate_normal(mean, cov, P)
    b = np.random.uniform(0.0, 2 * np.pi, P).reshape((P, 1))
    return G, b

### Post processing analysis ###
def training_loss(X_train, y_train, trainer, alpha=1e-7, reg=0.1):
    training_loss = []
    for num_iter in iterations:
        model = trainer(X_train, y_train, alpha, reg, num_iter)
        pred_y_train = predict(model, X_train)
        accuracy = metrics.accuracy_score(y_train, pred_y_train)
        print("Train accuracy with " + str(num_iter) + " iterations:
{0}".format(accuracy))
        training_loss.append(accuracy)
    return training_loss

def plot_training_loss(training_loss, title):
    plt.plot(iterations, training_loss, 'r-')
    plt.axis([0, 105000, .6, 1])
    plt.title(title)
    plt.show()

if __name__ == "__main__":
    X_train, X_test, y_train = load_dataset()

    # Uncomment the training_loss calls if you want to see how data was
    fitted

    """ Training losses for gradient descent """
    # print("Batch gradient descent on log lifted data")
    # log_training_loss = training_loss(log_lift(X_train), y_train, trai
n_gd)
    log = [0.60521739130434782, 0.93275362318840582,

```



```

0.94260869565217387, 0.94376811594202903, 0.94492753623188408, 0.94666666
6666666666, 0.94608695652173913, 0.94608695652173913, 0.9466666666666666
6, 0.94608695652173913]
# print("Batch gradient descent on standardized data")
# standardized_training_loss = training_loss(standardize(X_train), y
_train, train_gd)
standardized = [0.65884057971014498, 0.90927536231884054, 0.91130434
782608694, 0.9127536231884058, 0.91449275362318838, 0.91565217391304343,
0.9147826086956522, 0.91594202898550725, 0.91594202898550725, 0.9162318
8405797107]
# print("Batch gradient descent on binarized data")
# binarized_training_loss = training_loss(binarize(X_train), y_train,
train_gd)
binarized = [0.39478260869565218, 0.87971014492753619, 0.89246376811
594208, 0.89507246376811589, 0.90000000000000002, 0.90318840579710147,
0.90550724637681157, 0.90666666666666662, 0.91188405797101446, 0.9124637
6811594199]

""" Training losses for stochastic gradient descent """
# print("Stochastic gradient descent on log lifted data")
# sgd_log_training_loss = training_loss(log_lift(X_train), y_train,
train_sgd, alpha=1e-4)
sgd_log = [0.60521739130434782, 0.92695652173913046, 0.9237681159420
2901, 0.9301449275362319, 0.93710144927536232, 0.92985507246376808, 0.93
275362318840582, 0.92985507246376808, 0.93391304347826087, 0.93652173913
043479]
# print("Stochastic gradient descent on standardized data")
# sgd_standardized_training_loss = training_loss(standardize(X_train), y_train,
train_sgd, alpha=1e-4)
sgd_standardized = [0.65884057971014498, 0.9040579710144927, 0.90869
565217391302, 0.90927536231884054, 0.91130434782608694, 0.91130434782608
694, 0.9104347826086957, 0.91130434782608694, 0.91072463768115941, 0.911
30434782608694]
# print("Stochastic gradient descent on binarized data")
# sgd_binarized_training_loss = training_loss(binarize(X_train), y_train,
train_sgd, alpha=1e-3)
sgd_binarized = [0.39478260869565218, 0.89681159420289858, 0.8886956
52173913, 0.89217391304347826, 0.8863768115942029, 0.88579710144927537,
0.89478260869565218, 0.89101449275362321, 0.89130434782608692, 0.8811594
2028985506]

""" Training losses for stochastic gradient descent with decaying alpha """
# print("Decaying Alpha Stochastic gradient descent on log lifted data")
# sgda_log_training_loss = training_loss(log_lift(X_train), y_train,
train_sgd_decay, alpha=1e-3, reg=0.01)
# print("Training Losses: " + str(sgda_log_training_loss))
sgda_log = [0.60521739130434782, 0.60521739130434782, 0.605217391304
34782, 0.60521739130434782, 0.60521739130434782, 0.60521739130434782, 0.
60521739130434782, 0.60521739130434782, 0.60521739130434782, 0.605217391
30434782]
# print("Decaying Alpha Stochastic gradient descent on standardized
data")
# sgda_standardized_training_loss = training_loss(standardize(X_train), y_train,
train_sgd_decay, alpha=1e-3, reg=0.01)
# print("Training Losses: " + str(sgda_standardized_training_loss))

```

```

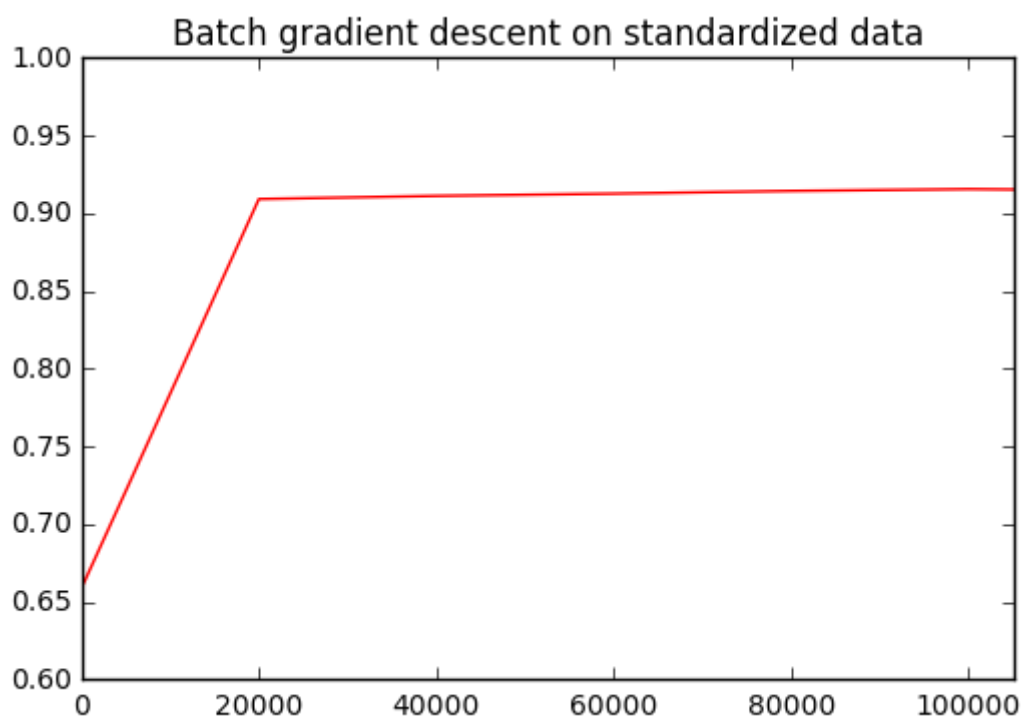
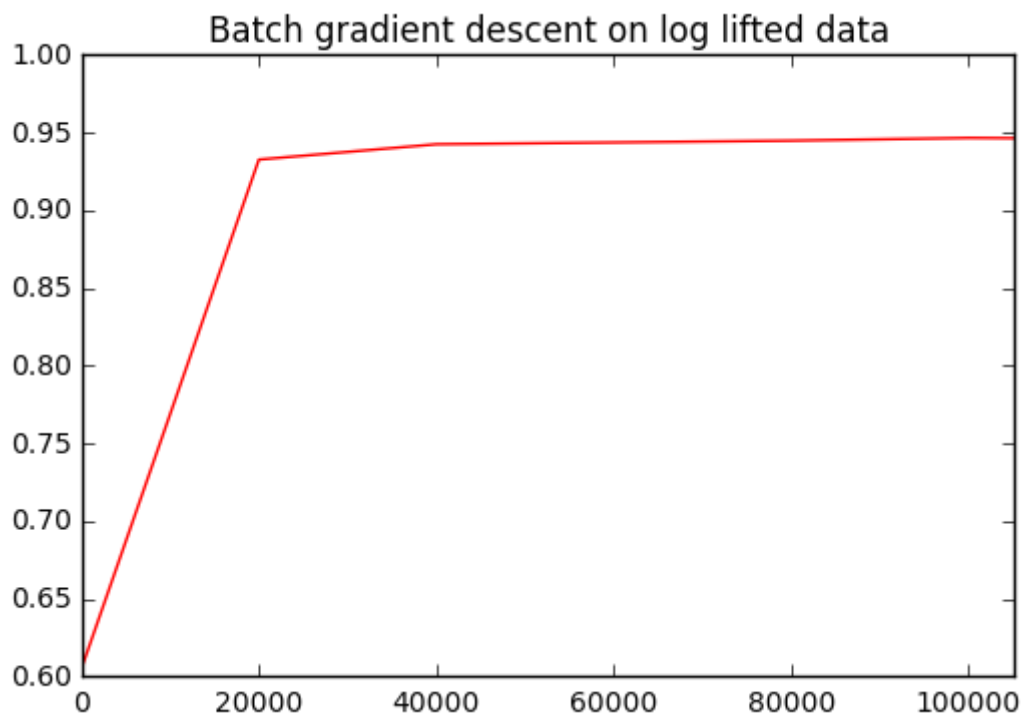
sgda_standardized = [0.65884057971014498, 0.88144927536231887, 0.892
75362318840579, 0.87449275362318846, 0.88144927536231887, 0.894202898550
72466, 0.87855072463768114, 0.888695652173913, 0.88550724637681155, 0.89
507246376811589]
# print("Decaying Alpha Stochastic gradient descent on binarized dat
a")
# sgda_binarized_training_loss = training_loss(binarize(X_train), y_
train, train_sgd_decay, alpha=1e-3, reg=0.01)
# print("Training Losses: " + str(sgda_binarized_training_loss))
sgda_binarized = [0.39478260869565218, 0.70028985507246377, 0.864057
97101449278, 0.87333333333333329, 0.60521739130434782, 0.878550724637681
14, 0.83043478260869563, 0.71739130434782605, 0.82521739130434779, 0.718
507246376811]

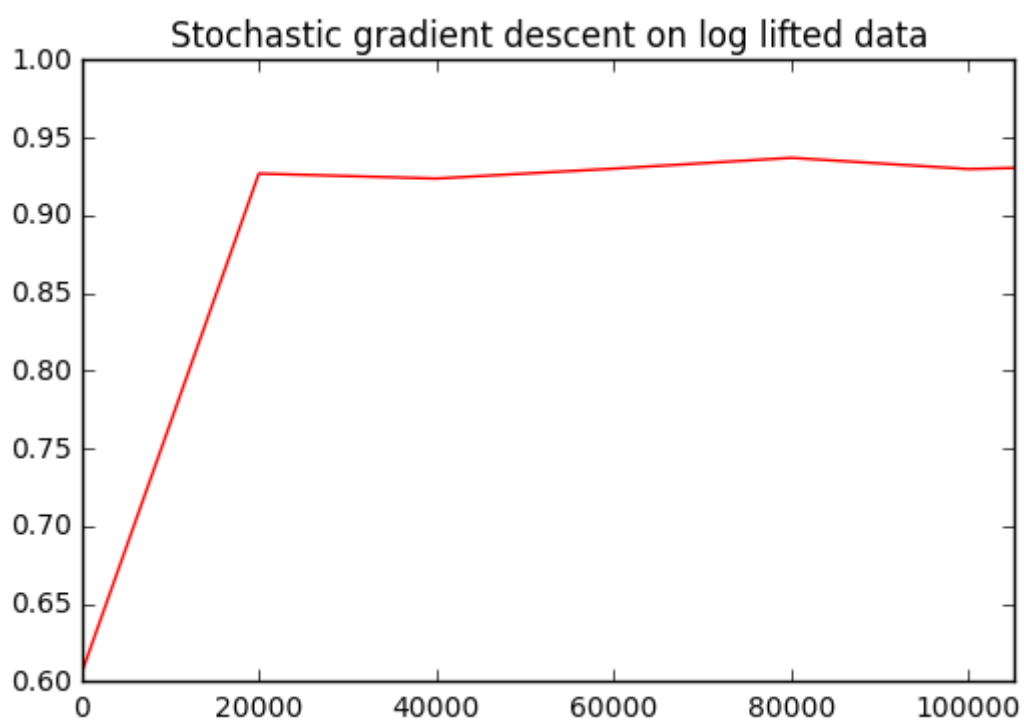
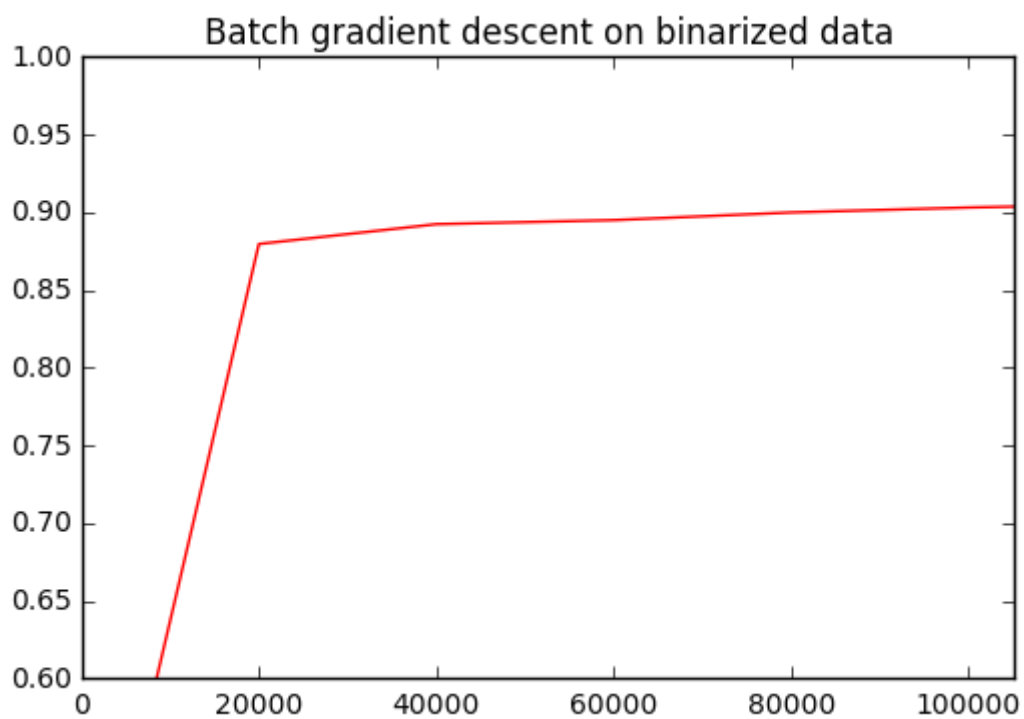
""" Plotting """
plot_training_loss(log, "Batch gradient descent on log lifted data")
plot_training_loss(standardized, "Batch gradient descent on standard
ized data")
plot_training_loss(binarized, "Batch gradient descent on binarized d
ata")
plot_training_loss(sgd_log, "Stochastic gradient descent on log lift
ed data")
plot_training_loss(sgd_standardized, "Stochastic gradient descent on
standardized data")
plot_training_loss(sgd_binarized, "Stochastic gradient descent on bi
narized data")
plot_training_loss(sgda_log, "Stochastic gradient descent with decay
ing alpha on log lifted data")
plot_training_loss(sgda_standardized, "Stochastic gradient descent w
ith decaying alpha on standardized data")
plot_training_loss(sgda_binarized, "Stochastic gradient descent with
decaying alpha on binarized data")

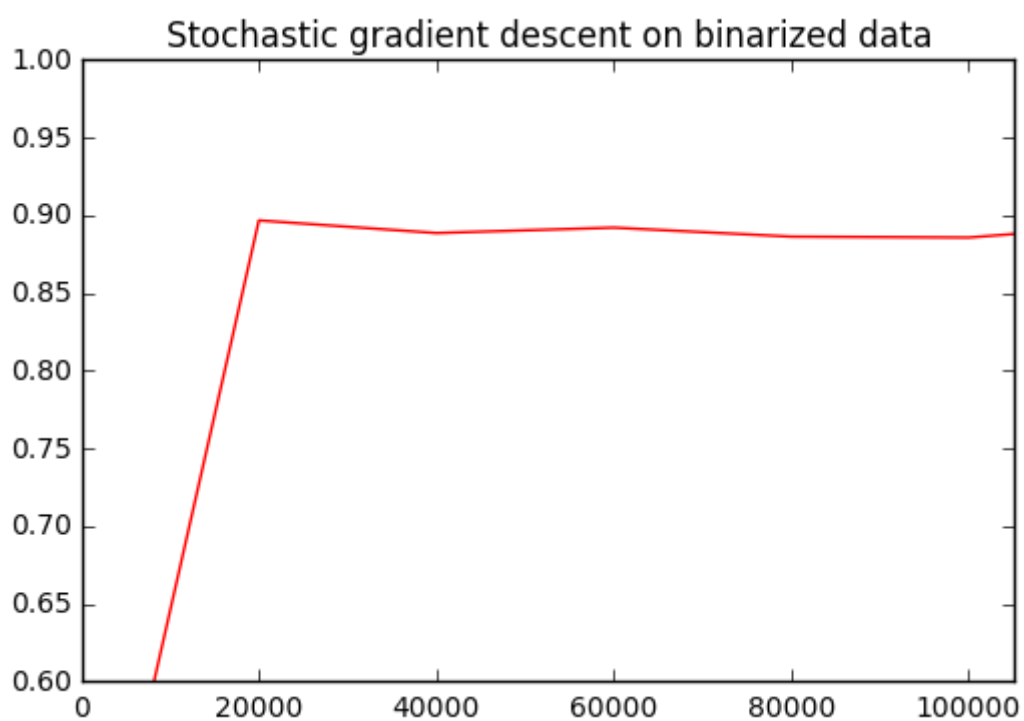
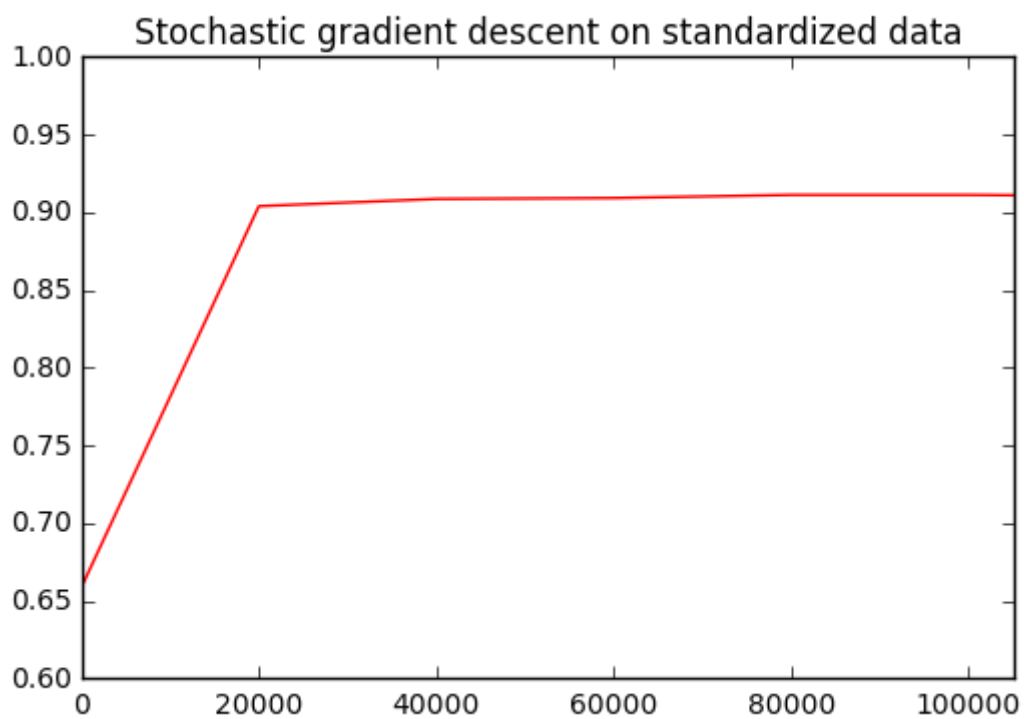
""" Kaggle Test """
model = train_gd(log_lift(X_train), y_train, alpha=1e-7, reg=0.1, nu
m_iter=100000)
pred_y_test = predict(model, log_lift(X_test))

c = csv.writer(open("hw3_kaggle.csv", "wt"))
c.writerow(['Id', 'Category'])
for i, val in enumerate(pred_y_test):
    c.writerow((i+1, int(val)))

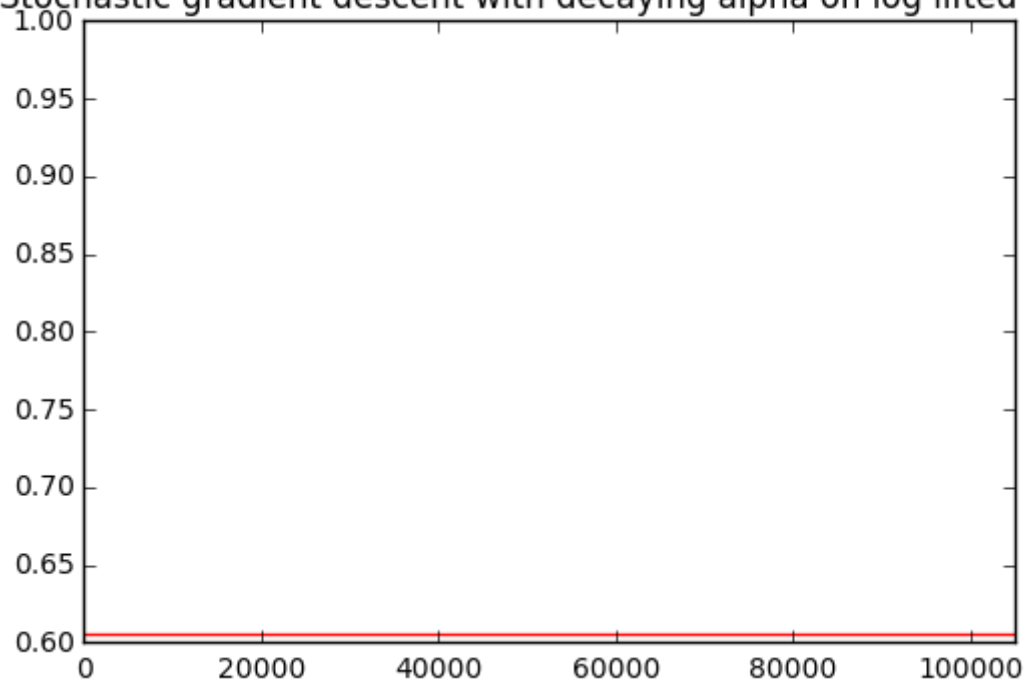
```



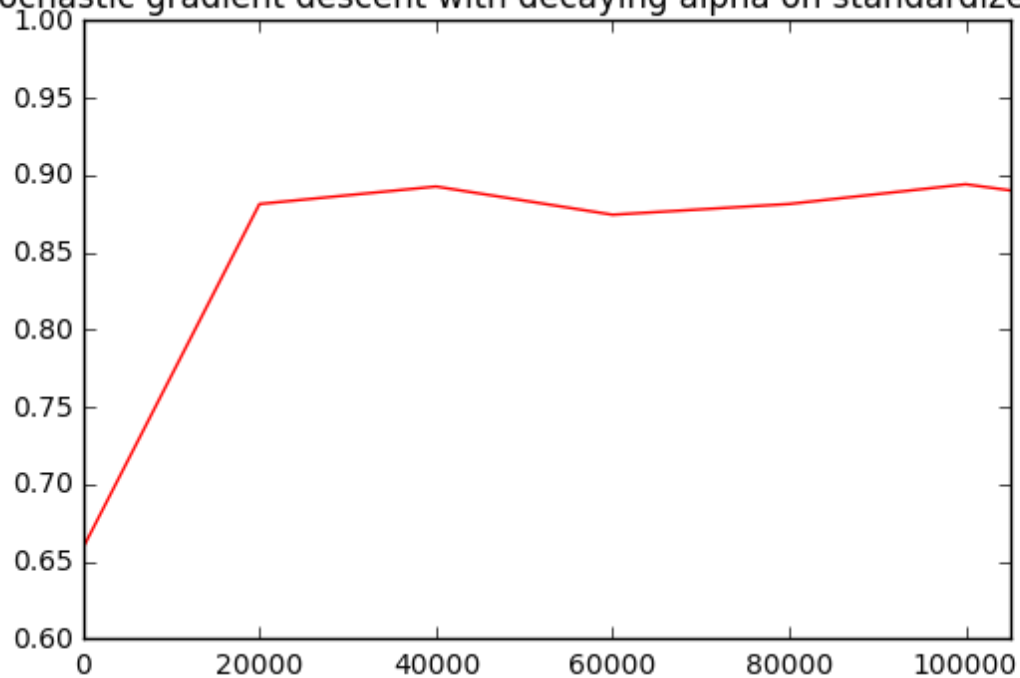




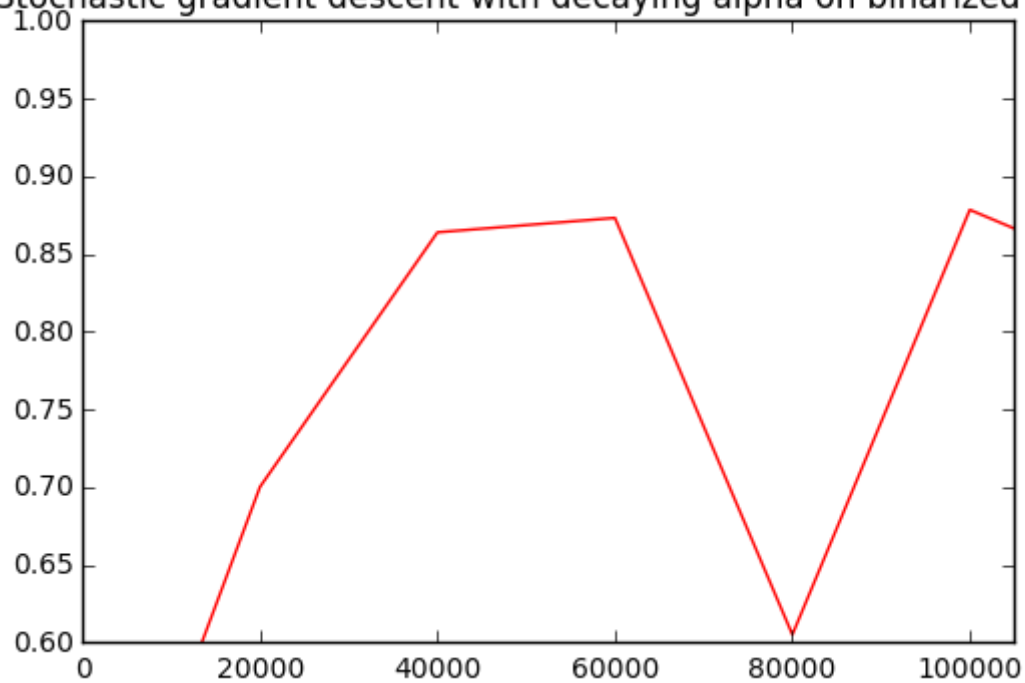
Stochastic gradient descent with decaying alpha on log lifted data



Stochastic gradient descent with decaying alpha on standardized data



Stochastic gradient descent with decaying alpha on binarized data



In []: