

This exam-prep discussion section covers regression.

## 1 Multiple Choice

(2) [3 pts] Duplicating a feature in linear regression

☐ Can reduce the L2-Penalized Residual Sum of Squares.

☐ Can reduce the L1-Penalized Residual Sum of Squares (RSS).

☐ Does not reduce the Residual Sum of Squares (RSS).

☐ None of the above

(24) [3 pts] Which of the following statements are true for a design matrix  $X \in \mathbb{R}^{n \times d}$  with  $d > n$ ? (The rows are  $n$  sample points and the columns represent  $d$  features.)

☐ Least-squares linear regression computes the weights  $w = (X^T X)^{-1} X^T y$ .

☐ The sample points are linearly separable.

☐  $X$  has exactly  $d - n$  eigenvectors with eigenvalue zero.

☐ At least one principal component direction is orthogonal to a hyperplane that contains all the sample points.

(n) [3 pts] Let  $w^*$  be the solution you obtain in standard least-squares linear regression. What solution do you obtain if you scale all the input features (but not the labels  $y$ ) by a factor of  $c$  before doing the regression?

☐  $\frac{1}{c} w^*$

☐  $c w^*$

☐  $\frac{1}{c^2} w^*$

☐  $c^2 w^*$

(o) [3 pts] In least-squares linear regression, adding a regularization term can

☐ increase training error.

☐ increase validation error.

☐ decrease training error.

☐ decrease validation error.

(p) [3 pts] You have a design matrix  $X \in \mathbb{R}^{n \times d}$  with  $d = 100,000$  features and vector  $\mathbf{y} \in \mathbb{R}^n$  of binary 0-1 labels. When you fit a logistic regression model to your design matrix, your test error is much worse than your training error. You suspect that many of the features are useless and are therefore causing overfitting. What are some ways to eliminate the useless features?

☐ Use  $\ell_1$  regularization.

☐ Use  $\ell_2$  regularization.

☐ Iterate over features; check if removing feature  $i$  increases validation error; remove it if not.

☐ If the  $i$ th eigenvalue  $\lambda_i$  of the sample covariance matrix is 0, remove the  $i$ th feature/column.

## 2 L2-Regularized Linear Regression with Newton's Method (Spring 2014)

Recall that the objective function for L2-regularized linear regression is

$$J(\mathbf{w}) = \|X\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

where  $X$  is the design matrix (the rows of  $X$  are the data points).

The global minimizer of  $J$  is given by:

$$\mathbf{w}^* = (X^T X + \lambda I)^{-1} X^T \mathbf{y}$$

(a) [8 pts] Consider running Newton's method to minimize  $J$ .

Let  $\mathbf{w}_0$  be an arbitrary initial guess for Newton's method. Show that  $\mathbf{w}_1$ , the value of the weights after one Newton step, is equal to  $\mathbf{w}^*$ .

### Q3. [10 pts] Error-Prone Sensors

We want to perform linear regression on the outputs of  $d$  building sensors measured at  $n$  different times, to predict the building's energy use. Unfortunately, some of the sensors are inaccurate and prone to large errors and, occasionally, complete failure. Fortunately, we have some knowledge of the relative accuracy and magnitudes of the sensors.

Let  $X$  be a  $n \times (d + 1)$  design matrix whose first  $d$  columns represent the sensor measurements and whose last column is all 1's. (Each sensor column has been normalized to have variance 1.) Let  $y$  be a vector of  $n$  target values, and let  $w$  be a vector of  $d + 1$  weights (the last being a bias term  $\alpha$ ). We decide to minimize the cost function

$$J(w) = \|Xw - y\|_1 + \lambda w^T D w,$$

where  $D$  is a diagonal matrix with diagonal elements  $D_{ii}$  (with  $D_{d+1,d+1} = 0$  so we don't penalize the bias term).

- (a) [2 pts] Why might we choose to minimize the  $\ell_1$ -norm  $\|Xw - y\|_1$  as opposed to the  $\ell_2$ -norm  $\|Xw - y\|_2$  in this scenario?
- (b) [2 pts] Why might we choose to minimize  $w^T D w$  as opposed to  $|w'|^2$ ? What could the values  $D_{ii}$  in  $D$  represent?
- (c) [6 pts] Derive the batch gradient descent rule to minimize our cost function. Hint: let  $p$  be a vector with components  $p_i = \text{sign}(X_i^T w - y_i)$ , and observe that  $\|Xw - y\|_1 = (Xw - y)^T p$ . For simplicity, assume that no  $X_i^T w - y_i$  is ever exactly zero.