**INSTRUCTIONS**

- You have 70 minutes to complete the exam.

- The exam is closed book, closed notes, closed computer, closed calculator, except for two 8.5" × 11" crib sheets of your own creation.

- Mark your answers **on the exam itself**. We will *not* grade answers written on scratch paper.

| | |
|---|---|
| Last name | |
| First name | |
| Student ID number | |
| CalCentral email (`_@berkeley.edu`) | |
| Exam room | |
| Name of the person to your left | |
| Name of the person to your right | |
| *All the work on this exam is my own.* **(please sign)** | |

**Terminology and Notation Reference:**

| | |
|---|---|
| $\exp(x)$ | $e^x$ |
| $\log(x)$ | $\log_e x$ |
| Linear regression model | $E[Y|X] = X^T\beta$ |
| Logistic (or sigmoid) function | $\sigma(t) = \frac{1}{1+\exp(-t)}$ |
| Logistic regression model | $P(Y = 1|X) = \sigma(X^T\beta)$ |
| Squared error loss | $L(y, \theta) = (y - \theta)^2$ |
| Absolute error loss | $L(y, \theta) = |y - \theta|$ |
| Cross-entropy loss | $L(y, \theta) = -y\log\theta - (1 - y)\log(1 - \theta)$ |
| Bias | $\mathrm{Bias}[\hat{\theta}, \theta] = E[\hat{\theta}] - \theta$ |
| Variance | $\mathrm{Var}[\hat{\theta}] = E[(\hat{\theta} - E[\hat{\theta}])^2]$ |
| Mean squared error | $\mathrm{MSE}[\hat{\theta}, \theta] = E[(\hat{\theta} - \theta)^2]$ |

1. **(8 points)   Feature Engineering**

For each dataset depicted below in a scatterplot, fill in the squares next to **all** of the letters for the vector-valued functions $f$ that would make it possible to choose a column vector $\beta$ such that $y_i = f(x_i)^T \beta$ for all $(x_i, y_i)$ pairs in the dataset. The input to each $f$ is a scalar $x$ shown on the horizontal axis, and the corresponding $y$ value is shown on the vertical axis.
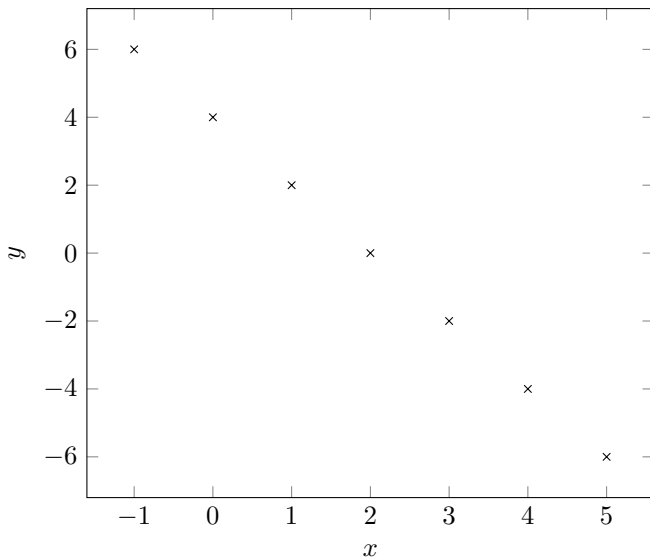
(A) $f(x) = [1 \quad x]^T$
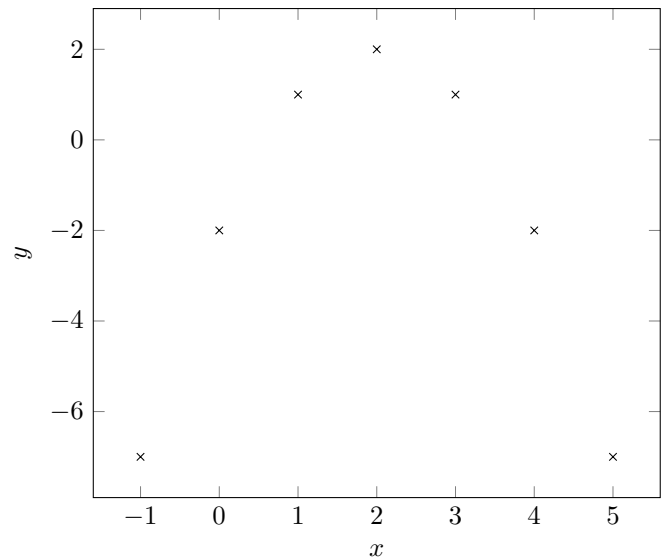
(B) $f(x) = [x \quad 2x]^T$

(C) $f(x) = [1 \quad x \quad x^2]^T$

(D) $f(x) = [1 \quad |x|]^T$

(E) None of the above
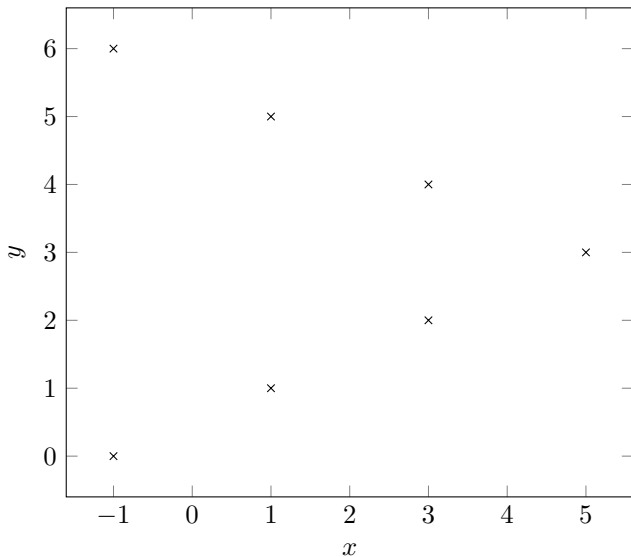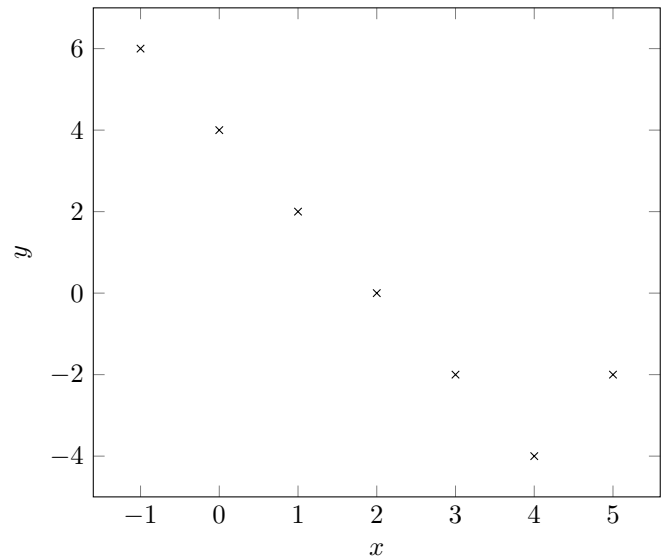
(i) (**2 pt**) ☐ A   ☐ B   ☐ C   ☐ D   ☐ E

(ii) (**2 pt**) ☐ A   ☐ B   ☐ C   ☐ D   ☐ E

(iii) (**2 pt**) ☐ A   ☐ B   ☐ C   ☐ D   ☐ E

(iv) (**2 pt**) ☐ A   ☐ B   ☐ C   ☐ D   ☐ E

*cross entropy loss*

## 2. (6 points)   Estimation

A learning set $(x_1, y_1), \ldots, (x_{10}, y_{10})$ is sampled from a population where $X$ and $Y$ are both binary.

The learning set data are summarized by the following table of row counts:

| $x$ | $y$ | Count |
|---|---|---|
| 0 | 0 | 2 |
| 0 | 1 | 3 |
| 1 | 0 | 1 |
| 1 | 1 | 4 |

$$L(y, \theta) = -y \log \theta - (1-y) \log(1-\theta)$$

**(a) (4 pt)** You decide to fit a constant model $P(Y = 1|X = 0) = P(Y = 1|X = 1) = \alpha$ using the cross-entropy loss function and no regularization. What is the formula for the empirical risk on this learning set for this model and loss function? What estimate of the model parameter $\alpha$ minimizes empirical risk? **You must show your work for finding the estimate $\hat{\alpha}$ to receive full credit.**

*Recall:* Since $Y$ is binary, $P(Y = 0|X) + P(Y = 1|X) = 1$ for any $X$.

$2 \times (0,0) :\quad L(y, \alpha) = -\log(1-\alpha)$

$1 \times (1,0): -\log(1-\alpha)$

Empirical Risk:   $3 \times (0,1) \quad L(y, \alpha) = -\log \alpha$

$4 \times (1,1): -\log \alpha$

Estimate $\hat{\alpha}$ (show your work):

$$R(\alpha) = \frac{-1}{10} \left( 7 \log \alpha + 3 \log(1-\alpha) \right)$$

$$\frac{\partial R}{\partial \alpha} = -\frac{1}{10} \left( \frac{7}{\alpha} + \frac{3}{1-\alpha}(-1) \right) = 0$$

$$\frac{7}{\alpha} = \frac{3}{1-\alpha}$$

$$7 - 7\alpha = 3\alpha$$

$$10\alpha = 7$$

$$\boxed{\hat{\alpha} = \frac{7}{10}}$$

**(b) (2 pt)** The true population probability $P(Y = 0|X = 0)$ is $\frac{1}{3}$. Provide an expression in terms of $\hat{\alpha}$ for the **bias** of the estimator of $P(Y = 0|X = 0)$ described in part (a) for the constant model. **You may use $E[\ldots]$ in your answer to denote an expectation under the data generating distribution of the learning set, but do not write $P(\ldots)$ in your answer.**

$$\hat{P}(y=0|X=0) = 1 - \alpha$$

$$\text{Bias}[\hat{P}(Y=0|X=0), P(Y=0|X=0)] = \underline{E[\hat{P}(y=0|X=0)] - P(y=0|X=0)}$$

$$= E[1-\alpha] - \frac{1}{3} = \frac{2}{3} - E[\alpha]$$

**3. (6 points)    Linear Regression**

A learning set of size four is sampled from a population where $X$ and $Y$ are both quantitative:

$$(x_1, y_1) = (2.5, 3)$$
$$(x_2, y_2) = (2, 5)$$
$$(x_3, y_3) = (1, 3)$$
$$(x_4, y_4) = (3, 5).$$

You fit a linear regression model $E[Y|X] = \beta_0 + X\beta_1$, where $\beta_0$ and $\beta_1$ are scalar parameters, by ridge regression, minimizing the following objective function:

$$\frac{1}{4}\left(\sum_{i=1}^{4} \overbrace{(y_i - (\beta_0 + x_i\beta_1))^2}^{\|y - X\beta\|_2^2} + 4\frac{\beta_0^2 + \beta_1^2}{3}\right).$$

(a) **(4 pt)** Fill in all blanks below to compute the parameter estimates that minimize this regularized empirical risk. (You do not need to compute their values; just fill in the matrices appropriately.)

$$\frac{1}{4}\sum_{i=1}^{4}(y_i - \hat{y}_i)^2 + \frac{1}{3}\|\beta\|_2^2$$
$$= \frac{1}{4}\left(\sum (y_i - \hat{y}_i)^2 + \frac{4}{3}\|\beta\|_2^2\right)$$

$$X_n^T = \begin{bmatrix} 1 & 1 & 1 & 1 \\ \hline 2.5 & 2 & 1 & 3 \end{bmatrix}$$

$$Y_n^T = \begin{bmatrix} 3 & 5 & 3 & 5 \end{bmatrix}$$

$$\begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix} = \left(X_n^T X_n + \begin{bmatrix} 4/3 & 0 \\ \hline 0 & 4/3 \end{bmatrix}\right)^{-1} X_n^T Y_n.$$

(b) **(2 pt)** Without computing values for $\hat{\beta}_0$ and $\hat{\beta}_1$, write an expression for the squared error loss of the learning set observation $(x_4, y_4)$ in terms of $\hat{\beta}_0$ and $\hat{\beta}_1$ and any relevant numbers. **Your solution should not contain any of $\hat{y}_4$, $x_4$, or $y_4$, but instead just numbers and $\hat{\beta}_0$ and $\hat{\beta}_1$.**

$$L(y_4, \hat{y}_4) = \left(5 - (\beta_0 + 3\beta_1)\right)^2 \qquad (x_4, y_4) = (3, 5)$$

**4. (8 points)  Model Selection**

(a) **(2 pt)** You have a quantitative outcome $Y$ and two quantitative covariates $(X_1, X_2)$. You want to fit a linear regression model for the conditional expected value $E[Y|X]$ of the outcome given the covariates, including an intercept. Bubble in the **minimum** dimension of the parameter vector $\beta$ needed to express this linear regression model?

○ 1    ○ 2    ○ 3    ○ 4    ○ 5    ○ 6    ○ 7    ○ None of these

(b) **(2 pt)** You have a quantitative outcome $Y$ and two qualitative covariates $(X_1, X_2)$. $X_1 \in \{a, b, c, d\}$, $X_2 \in \{e, f, g\}$, and there is no ordering to the values for either variable. You want to fit a linear regression model for the conditional expected value $E[Y|X]$ of the outcome given the covariates, including an intercept. Bubble in the **minimum** dimension of the parameter vector $\beta$ needed to express this linear regression model?

○ 2   ○ 3   ○ 4   ○ 5   ○ 6   ○ 7   ○ 8   ○ 9   ○ 10   ○ 11   ○ 12   ○ 13

(c) **(2 pt)** Bubble all true statements: In ridge regression, when the assumptions of the linear model are satisfied, the larger the shrinkage/penalty parameter,

☐ the larger the magnitude of the bias of the estimator of the regression coefficients $\beta$.

☐ the smaller the magnitude of the bias of the estimator of the regression coefficients $\beta$.

☐ the larger the variance of the estimator of the regression coefficients $\beta$.

☐ the smaller variance of the estimator of the regression coefficients $\beta$.

☐ the smaller the true mean squared error of the estimator of the regression coefficients $\beta$.

(d) **(2 pt)** Bubble all true statements: A good approach for selecting the shrinkage/penalty parameter in LASSO is to:

☐ minimize the learning set risk for the squared error ($L_2$) loss function.

☐ minimize the learning set risk for the absolute error ($L_1$) loss function.

☐ minimize the cross-validated regularized risk for the squared error ($L_2$) loss function.

☐ minimize the cross-validated risk for the squared error ($L_2$) loss function.

☐ minimize the variance of the estimator of the regression coefficients.

$$P(y=1|x) = \frac{1}{1+e^{-x^T\beta}} \qquad P(y=0|x) = \frac{e^{-x^T\beta}}{1+e^{-x^T\beta}}$$

$$\Rightarrow \frac{P(y=1|x)}{P(y=0|x)} = \frac{\frac{1}{1+e^{-x^T\beta}}}{\frac{e^{-x^T\beta}}{1+e^{-x^T\beta}}} = \frac{1}{e^{-x^T\beta}} = e^{x^T\beta}$$

**5. (12 points)   Logistic Regression**

(a) **(2 pt)** Bubble the expression that describes the odds ratio $\frac{P(Y=1|X)}{P(Y=0|X)}$ of a logistic regression model.
*Recall*: $P(Y=0|X) + P(Y=1|X) = 1$ for any $X$.

○ $X^T\beta$      ○ $-X^T\beta$      ● $\exp(X^T\beta)$      ○ $\sigma(X^T\beta)$      ○ None of these

(b) **(2 pt)** Bubble the expression that describes $P(Y=0|X)$ for a logistic regression model.

● $\sigma(-X^T\beta)$      ○ $1-\log(1+\exp(X^T\beta))$      ○ $1+\log(1+\exp(-X^T\beta))$      ○ None of these

(c) **(2 pt)** Bubble **all** of the following that are typical effects of adding an $L_1$ regularization penalty to the loss function when fitting a logistic regression model with parameter vector $\beta$.
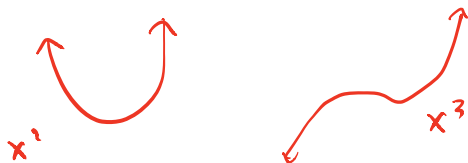
☐ The magnitude of the elements of the estimator of $\beta$ are increased.

☑ The magnitude of the elements of the estimator of $\beta$ are decreased.

☐ All elements of the estimator of $\beta$ are non-negative.

☑ Some elements of the estimator of $\beta$ are zero.

☐ None of the above.

(d) **(3 pt)** What would be the primary disadvantage of a regularization term of the form $\sum_{j=1}^{J}\beta_j^3$ rather than the more typical ridge penalty $\sum_{j=1}^{J}\beta_j^2$ for logistic regression? Answer in one sentence.

$\sum \beta_j^3$ has no global min,

so $\beta_i \to \infty$

$x^2$  $x^3$

(e) **(3 pt)** For a logistic regression model $P(Y=1|X) = \sigma(-2-3X)$, where $X$ is a scalar random variable, what values of $x$ would give $P(Y=0|X=x) \geq \frac{3}{4}$? **You must show your work for full credit.**

$$P(y=0|x) = 1 - P(y=1|x) = 1 - \sigma(-2-3X)$$

$$1 - \sigma(-2-3x) \geq \frac{3}{4}$$

$$\sigma(-2-3x) \leq \frac{1}{4}$$

$$\frac{1}{1+e^{-(-2-3x)}} = \frac{1}{1+e^{2+3x}} \leq \frac{1}{4}$$

$$1+e^{2+3x} \geq 4$$

$$e^{2+3x} \geq 3$$

$$2 + 3x \geq \log 3$$

$$x \geq \frac{\log 3 - 2}{3}$$