# Ordinary Least Squares

Multiple Linear Regression Model: $\hat{\mathbb{Y}} = \mathbb{X}\theta$ with design matrix $\mathbb{X}$, response vector $\mathbb{Y}$, and predicted vector $\hat{\mathbb{Y}}$. If there are $p$ features plus a bias/intercept, then the vector of parameters $\theta = [\theta_0, \theta_1, \ldots, \theta_p]^T \in \mathbb{R}^{p+1}$. The vector of estimates $\hat{\theta}$ is obtained from fitting the model to the sample $(\mathbb{X}, \mathbb{Y})$.

| Concept | Formula | Concept | Formula |
|---|---|---|---|
| Mean squared error | $R(\theta) = \frac{1}{n}\|\mathbb{Y} - \mathbb{X}\theta\|_2^2$ | Normal equation | $\mathbb{X}^T\mathbb{X}\hat{\theta} = \mathbb{X}^T\mathbb{Y}$ |
| Least squares estimate, if $\mathbb{X}$ is full rank | $\hat{\theta} = (\mathbb{X}^T\mathbb{X})^{-1}\mathbb{X}^T\mathbb{Y}$ | Residual vector, $e$ | $e = \mathbb{Y} - \hat{\mathbb{Y}}$ |
| | | Multiple $R^2$ (coefficient of determination) | $R^2 = \dfrac{\text{variance of fitted values}}{\text{variance of } y}$ |

# Regularization

| Concept | Formula | Concept | Formula |
|---|---|---|---|
| Ridge Regression L2 Regularization | $\frac{1}{n}\|\mathbb{Y} - \mathbb{X}\theta\|_2^2 + \lambda\|\theta\|_2^2$ | Squared L2 Norm of $\theta \in \mathbb{R}^p$ | $\|\theta\|_2^2 = \sum_{j=1}^p \theta_j^2$ |
| Ridge regression estimate (closed form) | $\hat{\theta}_{\text{ridge}} = (\mathbb{X}^T\mathbb{X} + n\lambda I)^{-1}\mathbb{X}^T\mathbb{Y}$ | | |
| LASSO Regression L1 Regularization | $\frac{1}{n}\|\mathbb{Y} - \mathbb{X}\theta\|_2^2 + \lambda\|\theta\|_1$ | L1 Norm of $\theta \in \mathbb{R}^p$ | $\|\theta\|_1 = \sum_{j=1}^p |\theta_j|$ |

# Gradient Descent

Let $L$ be an objective function to minimize with respect to $\theta$; assume that some optimal parameter vector $\hat{\theta}$ exists. Suppose $\theta^{(0)}$ is some starting estimate at $t = 0$, and $\theta^{(t)}$ is the estimate at step $t$. Then for a learning rate $\alpha$, the gradient update step to compute $\theta^{(t+1)}$ is:

$$\theta^{(t+1)} = \theta^{(t)} - \alpha\nabla_\theta L$$

where $\nabla_\theta L$ is the partial derivative/gradient of $L$ with respect to $\theta$, evaluated at $\theta^{(t)}$.

# Classification and Logistic Regression

### Confusion Matrix

Columns are the predicted values $\hat{y}$ and rows are the actual classes $y$.

| | $\hat{y} = 0$ | $\hat{y} = 1$ |
|---|---|---|
| $y = 0$ | True negative (TN) | False Positive (FP) |
| $y = 1$ | False negative (FN) | True Positive (TP) |

### Classification Performance

Suppose you predict $n$ datapoints.

| Metric | Formula | Other Names |
|---|---|---|
| Accuracy | $\frac{TP+TN}{n}$ | |
| Precision | $\frac{TP}{TP+FP}$ | |
| Recall/TPR | $\frac{TP}{TP+FN}$ | True Positive Rate, Sensitivity |
| FPR | $\frac{FP}{FP+TN}$ | False Positive Rate, Specificity |

An ROC curve visualizes TPR vs. FPR for different thresholds $T$.

**Logistic Regression Model**: For input feature vector $x$, $\hat{P}_\theta(Y = 1|x) = \sigma(x^T\theta)$, where $\sigma(z) = 1/(1 + e^{-z})$. The estimate $\hat{\theta}$ is the parameter $\theta$ that minimizes the average cross-entropy loss on training data. For a single datapoint, define cross-entropy loss as $-[y\log(p) + (1 - y)\log(1 - p)]$, where $p$ is the probability that the response is 1.

**Logistic Regression Classifier**: For a given input $x$ and trained logistic regression model with parameter $\theta$, compute $p = \hat{P}(Y = 1|x) = \sigma(x^T\theta)$. predict response $\hat{y}$ with classification threshold $T$ as follows:

$$\hat{y} = \text{classify}(x) = \begin{cases} 1 & p \geq T \\ 0 & \text{otherwise} \end{cases}$$