

Linear Regression (9/11/13)

Lecturer: Barbara Engelhardt Scribes: Zachary Abzug, Mike Gloude-mans, Zhuosheng Gu, Zhao Song

1 Why use linear regression?

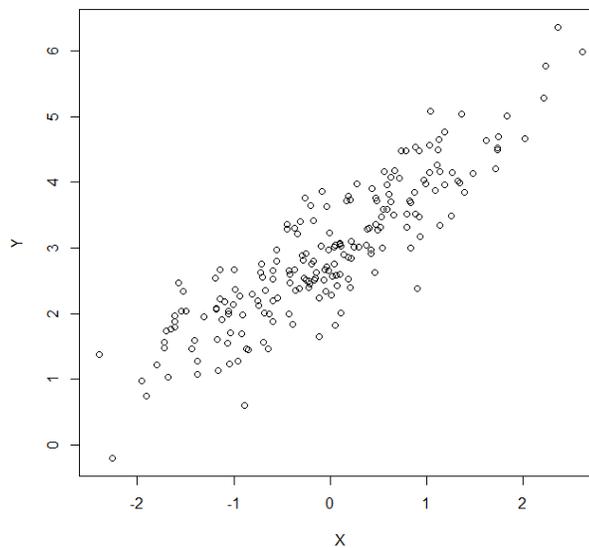


Figure 1: Scatter plot showing the relation between two variables, x and y

The general idea behind linear models is to capture how two observed variables are related in a linear way, assuming Gaussian noise. Given a fitted linear model (i.e., we have estimated the model parameters from training data), we may use this linear regression to:

- predict the value for y^* given a new observation x^*
- test if there is a linear relationship between two variables x and y
- when y is binary or multinomial, we can use this model to classify a new observation x^* .

Later in the course, we will learn non-linear and non-Gaussian versions of this basic model. Let's think about a bunch of possible relationships we want to model:

- x = age and y = height
- x = distance from the shore and y = weight of clams found

- $x =$ gestational age and $y =$ birthweight

Once we have fitted a linear model, then, given a new observation x^* , we can predict the corresponding value y^* . This model can also be used to test for associations between two variables. If there is no association between our variables, then the distribution of y should not depend conditionally on the value of x .

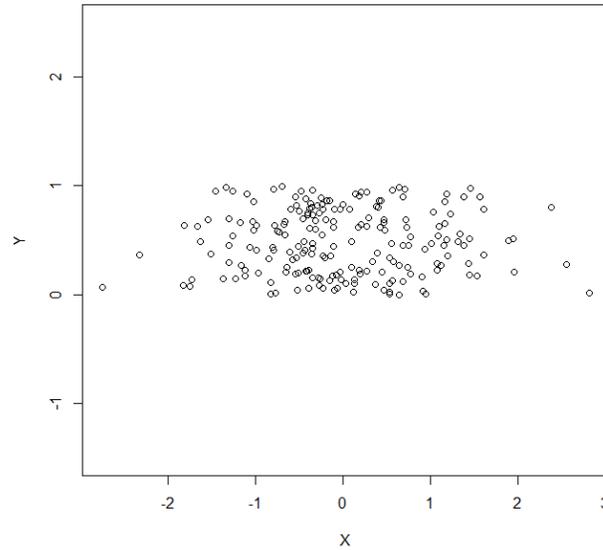


Figure 2: Scatter plot showing no relationship between two variables, x and y

2 Definitions & model specification

$y \in \mathbb{R}$: response

$\mathbf{x} \in \mathbb{R}^p$: predictors (also called covariates or explanatory variables)

$\boldsymbol{\beta} \in \mathbb{R}^p$: regression coefficients (also called effects)

Our training set consists of observations for n samples (\mathbf{x}, y) . To fit the model, we will estimate the coefficients. Note that bold font for a random variable represents a vector.

2.1 Linear model

Our linear model has the following form:

$$y = \mathbf{x}^T \boldsymbol{\beta} + \epsilon \quad (1)$$

where $\epsilon \sim \mathcal{N}(0, \sigma^2)$. Consequently, we have $y \mid \mathbf{x}^T, \boldsymbol{\beta}, \sigma^2 \sim \mathcal{N}(\mathbf{x}^T \boldsymbol{\beta}, \sigma^2)$. Note here that the distribution of \mathbf{x} doesn't *really* matter, as long as y is still distributed as a normal conditional on \mathbf{x} .

To build our linear model, we must first define our data. Let our training data $\mathcal{D}_{train} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$, where each \mathbf{x}_i is a $p \times 1$ column vector of p different predictors, and we have n observations of each predictor and response. Then, using our knowledge of multivariate normals, we can expand (1) to

$$\begin{bmatrix} y_1 \\ \cdot \\ \cdot \\ \cdot \\ y_n \end{bmatrix} = \mathcal{N}_n \left(\begin{bmatrix} \mathbf{x}_1^T \\ \cdot \\ \cdot \\ \cdot \\ \mathbf{x}_n^T \end{bmatrix} \beta, \text{diag}_n(\sigma^2) \right)$$

where X is an $n \times p$ matrix. Our covariance matrix is a diagonal matrix with a single variance term for all samples (σ^2), and all covariance terms are zero, indicating independence across the n dimensions of this Gaussian (i.e., over the samples). When our model is fitted, if we are given a new \mathbf{x}^* , we can predict the associated y^* value; in particular, $E[y^* | \mathbf{x}^*] = \mathbf{x}^{*T} \beta$. Furthermore, if we take many samples with the same x -value, we should expect to see a normal distribution of y -values, centered around the mean that we've predicted.

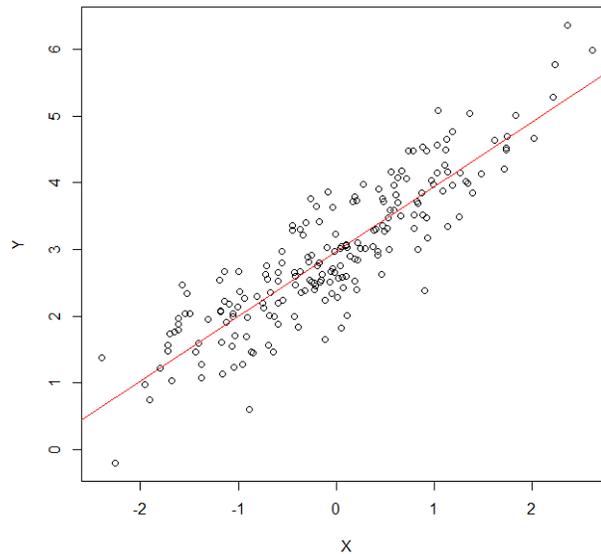


Figure 3: Our linear regression fit to the data. For any new x^* , we can use the model to predict the associated y^* by projecting the point \mathbf{x}^* onto the red regression line and identifying the corresponding value of y^* .

2.2 What about offsets?

Our regression line rarely will go through the origin (a zero year old person is not zero inches tall), so we must also include an intercept term, β_0 . We accomplish this by adding a 1 as the first element of our x vector, and a corresponding β_0 term as the first element of our β vector:

$$\mathbf{x} = [1, x_1, x_2, \dots, x_p]^T$$

Our model equation will then take the form $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$, where we can read off the value of the intercept from β_0 .

3 Evaluating our model

3.1 Residuals

How well do our response values we predict using the fitted model differ from the actual response values? We provide the following notation first (univariate case):

$$\begin{aligned} r_i &= y_i - \hat{\beta}x_i = y_i - \hat{y}_i \\ \hat{\beta} &= \text{our trained } \beta \\ \hat{y}_i &= \text{predicted value} \end{aligned}$$

These terms r_i are called the *residuals*. For each data point, the error is the same as the residual: $\epsilon_i = y_i - \hat{\beta}x_i = r_i$. Additionally, by our previous definition of the error term, $\mathbb{E}[r_i] = \mathbb{E}[\epsilon_i] = \mathbb{E}[y_i - \hat{\beta}x_i] = 0$. In other words, the expected residual is zero according to the Gaussian model.

3.2 Residual Sum of Squares (and friends!)

There are a number of metrics to quantify the fit of the model to the test set or training set data based on the residual values:

- Residual Sum of Squares (RSS): $\text{RSS}(\hat{\beta}, \mathcal{D}) = \sum_{i=1}^n (y_i - \hat{\beta}^T x_i)^2$
- Mean Squared Error (MSE): $\text{MSE} = \frac{1}{n} \text{RSS}(\hat{\beta}, \mathcal{D})$
- Root-Mean-Square Error (RMSE): $\text{RMSE} = \sqrt{\text{MSE}}$

4 Parameter Estimation

Our data for linear regression is represented as: $\mathcal{D} = \{(x_1, y_1), \dots, (x_n, y_n)\}$. For a univariate model ($p = 1$), where each x_i is a scalar value, the maximum likelihood estimate (MLE) for β , $\hat{\beta}$, can be found by taking the log likelihood, differentiating with respect to β , and setting equal to 0 and solving for β .

$$\begin{aligned} \ell(\beta; \mathcal{D}) &= \sum_{i=1}^n \log \left[\left(\frac{1}{2\pi\sigma^2} \right)^{1/2} \exp \left\{ \frac{-1}{2\sigma^2} (y_i - \beta x_i)^2 \right\} \right] \\ &= -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta x_i)^2 - \frac{n}{2} \log(2\pi\sigma^2) \\ &= -\frac{1}{2\sigma^2} \text{RSS}(\beta, \mathcal{D}) + c \end{aligned}$$

where c is not a function of β , and may be ignored because it is constant with respect to β . By looking at this equation, we see that we can maximize the log likelihood by equivalently minimizing the residual sum of squares (RSS).

Let's return to the general multivariate case, i.e. $p > 1$, so each data point x_i is a p -dimensional vector. We derive the MLE estimate for β as follows:

$$\begin{aligned} \frac{\partial [-\ell(\beta; \mathcal{D})]}{\partial \beta} &= \frac{\partial [\frac{1}{2}(\mathbf{y} - X\beta)^T(\mathbf{y} - X\beta)]}{\partial \beta} \\ &= X^T X \beta - X^T \mathbf{y} \end{aligned} \quad (2)$$

Setting (2) to zero, and solving for β , we have

$$\hat{\beta}_{MLE} = (X^T X)^{-1} X^T \mathbf{y} \quad (3)$$

Equation (3) is called the *normal equation*. This is a wonderful result: If $X^T X$ is invertible, then we have an efficient way of computing the MLE for the regression coefficients (the cost of a matrix inversion). However, if $X^T X$ is non-invertible (singular), then we cannot compute an estimate of β using the normal equation. One example of a singular matrix $X^T X$ happens when the number of samples n is smaller than the number of predictors p .

4.1 Ridge regression

To regularize this problem, and guard ourselves against singular matrices, we can include a prior distribution for β . One advantage of selecting a prior distribution for the regression coefficients is that it helps prevent distortions of the fit due to outliers. If we select a normal prior, the variance τ^2 of the coefficients reflects our confidence in the regression: a low value of τ^2 shrinks the coefficients more heavily toward zero, whereas a high value of τ^2 lets the data guide the selection of β more heavily. With

$$p(\beta | \tau^2) = \prod_{j=1}^p \mathcal{N}(\beta_j | 0, \tau^2)$$

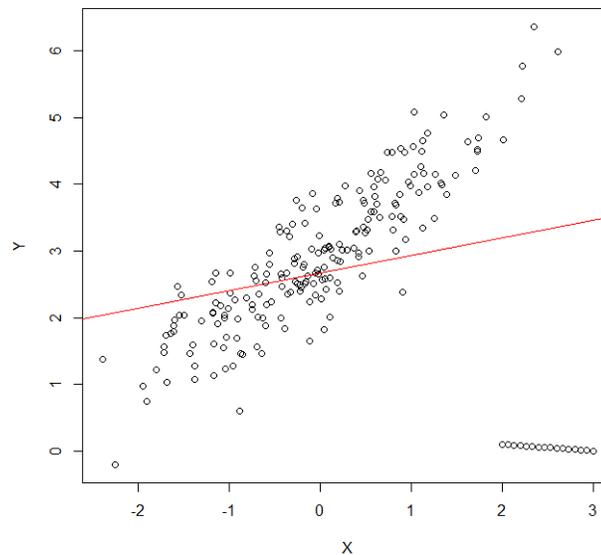


Figure 4: Distortion of the linear regression fit due to outliers.

we now calculate the log posterior distribution using Bayes Rule:

$$\begin{aligned} p(\boldsymbol{\beta}|\mathbf{x}, \mathbf{y}, \tau^2, \sigma^2) &\propto \sum_{i=1}^n \log \left[\left(\frac{1}{2\pi\sigma^2} \right)^{1/2} \exp \left\{ \frac{-1}{2\sigma^2} (y_i - \boldsymbol{\beta}^T \mathbf{x}_i)^2 \right\} \right] + \sum_{j=1}^p \log \mathcal{N}(\beta_j|0, \tau^2) \\ &= \frac{1}{n} \sum_{i=1}^n (y_i - \boldsymbol{\beta}^T \mathbf{x}_i)^2 + \frac{\sigma^2}{\tau^2} \sum_{j=1}^p \beta_j^2 \end{aligned} \quad (4)$$

The solution for $\hat{\boldsymbol{\beta}}_{MAP}$ in Eqn. (4) is referred to as *ridge regression*

$$\hat{\boldsymbol{\beta}}_{MAP} = \underset{\boldsymbol{\beta}}{\operatorname{argmax}} \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^p \beta_j^2$$

where the first term of the argmax is our MSE equation (see section 4.2) and the second term is referred to as the *penalty*. Here, λ can be regarded as a regularization parameter corresponding to $\frac{\sigma^2}{\tau^2}$ in Eq (4). We compute the MAP estimate for $\boldsymbol{\beta}$ using the same strategy as in the normal equation and obtain

$$\hat{\boldsymbol{\beta}}_{MAP} = (X^T X + \lambda I_p)^{-1} X^T \mathbf{y}$$

Notice here that $(X^T X + \lambda I_p)^{-1}$ exists as the matrix $(X^T X + \lambda I_p)$ is generally invertible (i.e., non-singular). In general, a simple way to make a singular matrix invertible is to add a very small value to the diagonal, in order to (artificially) create a full rank (and hence invertible) matrix.

5 Least mean squares (LMS) algorithm

Another way to solve for β is directly minimize the least square cost function (the residual sum of squares):

$$\begin{aligned} C(\beta) &= \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\ \frac{\partial C(\beta)}{\partial \beta} &= - \sum_{i=1}^n (y_i - \hat{y}_i) x_i \end{aligned}$$

Then we can use the steepest descent algorithm to find the estimate for β

$$\beta^{(t+1)} = \beta^{(t)} + \rho \sum_{i=1}^n (y_i - \beta^{(t)T} x_i) x_i$$

where ρ is the step size. Note here the summation is from 1 to n , which means we have to iterate through all of the n data instances before we update $\beta^{(t+1)}$. This is called a *batch method*. However, it is often tempting to update our estimate as we collect each sample: In some cases n is too large for machines to process at once, or the matrix inversion in the normal equation may be difficult if there are too many features p . So we would like a way to process the data one sample at a time, which is called an *online method*. The LMS algorithm can be rewritten as:

$$\beta^{(t+1)} = \beta^{(t)} + \rho (y_i - \beta^{(t)T} x_i) x_i,$$

which is known as *stochastic gradient descent*. This means at each iteration t , we randomly pick a data sample x_i, y_i from the data set and update our estimate of β until we convergence. This algorithm considers a single data point at a time, and is stochastic in selecting a sample from the full data.

5.1 Accounting for non-linearities

Imagine a model:

$$y_i = \phi(x_i)\boldsymbol{\beta} + \epsilon$$
$$\phi(x) = [1, x, x^2, \dots, x^p]$$

This would allow us to find non-linear relationships between x and y . We can also try combinations of x -values, for $x_i = [x_1, x_2]_i$:

$$\phi(x) = [x_1, x_2, x_1x_2, x_1^2x_2, \dots]$$

This function $\phi(x)$ is called the **basis function**. It turns out that the regression model above works with any basis function. This is still a *linear model* because the relationship between the coefficients β and y is linear: regardless of the $\phi(\cdot)$ term, this is a *linear model* because the latent parameter β enters the model in a linear way. Furthermore, the log likelihood is convex with respect to β , and this is true for all $\phi(x)$.

This nonlinear basis function will come up again when we talk about kernels, kernel methods, Gaussian processes, and adaptive basis methods.