



LINEAR REGRESSION

Buckle up...



Data Matrix

N the number of samples

p the number of features

$$\mathbf{X} \in \mathbb{R}^{N \times p}$$

$$\mathbf{X} = \begin{pmatrix} x_{00} & \cdots & x_{0p} \\ \vdots & \ddots & \vdots \\ x_{N0} & \cdots & x_{Np} \end{pmatrix}$$

Linear Model

We want to find some linear function which takes a single observation

$$X^T = (X_0, X_1, \dots, X_p)$$

and makes a prediction using a function of the form

$$f(X) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p$$

$$f(X) = \beta_0 + \sum_{j=1}^p X_j \beta_j$$

Linear Model

This can be written as a vector product:

$$f(X) = \beta_0 + \sum_{j=1}^p X_j \beta_j$$

is the same as

$$f(X) = \beta_0 + X^T \cdot \beta$$

Linear Regression

We wish to minimize the residual sum of squares:

$$RSS = \sum_{i=1}^N (y_i - f(x_i))^2$$

$$RSS = \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^p X_j \beta_j \right)^2$$

Linear Regression

Why do we choose to minimize this?

1. This is a reasonable measure of distance:
 - The quantity $(y_i - f(x_i))^2$ is always positive, as required for a metric
 - The distance from the set of true observations might as well be the sum of all the distances (ie. the distance from each prediction to the true observation, summed)
2. If we assume that our experimental errors are Gaussian then minimizing the RSS is equivalent to finding the most likely model under this assumption.
 - We will discuss this more later

Linear Regression

How to find the least-squares solution?

Write the RSS as a matrix-vector equation:

$$RSS = \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2$$

$$RSS = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

Linear Regression

To minimize RSS, take the usual approach of minimization:

$$\begin{aligned}\frac{\partial RSS}{\partial \beta} &= \frac{\partial \{(\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta)\}}{\partial \beta} \quad \star \\ 0 &= -2 \mathbf{X}^T (\mathbf{y} - \mathbf{X}\beta) \\ 0 &= \mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{X} \beta \\ \mathbf{X}^T \mathbf{X} \beta &= \mathbf{X}^T \mathbf{y}\end{aligned}$$

$$\begin{aligned}\hat{\beta} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ \hat{\mathbf{y}} &= \mathbf{X} \hat{\beta} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}$$

Some Matrix Rules

$$(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T$$

$$\frac{\partial \mathbf{x}^T \mathbf{a}}{\partial \mathbf{x}} = \frac{\partial \mathbf{a}^T \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a}$$

$$\frac{\partial \mathbf{a}^T \mathbf{X} \mathbf{b}}{\partial \mathbf{X}} = \mathbf{a} \mathbf{b}^T$$

$$\frac{\partial \mathbf{a}^T \mathbf{X}^T \mathbf{b}}{\partial \mathbf{X}} = \mathbf{b} \mathbf{a}^T$$

$$\frac{\partial \mathbf{a}^T \mathbf{X} \mathbf{a}}{\partial \mathbf{X}} = \frac{\partial \mathbf{a}^T \mathbf{X}^T \mathbf{a}}{\partial \mathbf{X}} = \mathbf{a} \mathbf{a}^T$$

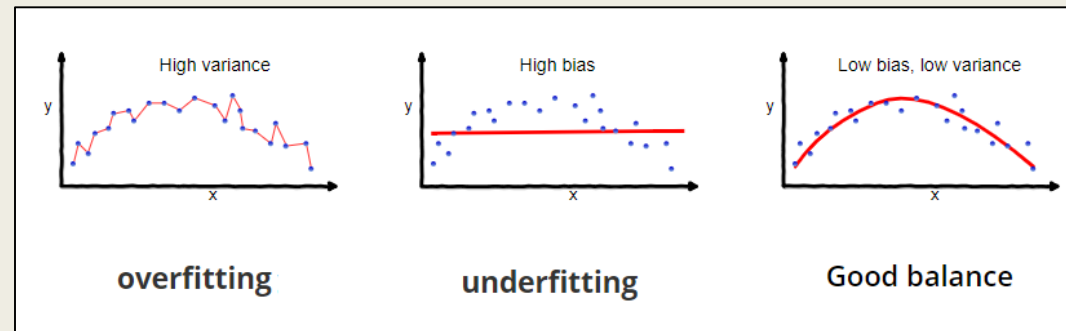
$$\frac{\partial \mathbf{b}^T \mathbf{X}^T \mathbf{X} \mathbf{c}}{\partial \mathbf{X}} = \mathbf{X} (\mathbf{b} \mathbf{c}^T + \mathbf{c} \mathbf{b}^T)$$

$$\frac{\partial}{\partial \mathbf{A}} (\mathbf{x} - \mathbf{A} \mathbf{s})^T \mathbf{W} (\mathbf{x} - \mathbf{A} \mathbf{s}) \quad \star \quad = -2 \mathbf{W} (\mathbf{x} - \mathbf{A} \mathbf{s}) \mathbf{s}^T$$

Linear Regression

Comments on least squares:

- The least squares solution is unique only if X is of full rank (all columns linearly independent)
- Gauss-Markov theorem: the least squares estimate is the minimum variance estimate amongst all unbiased linear estimates



- Least squares solutions often have high variance but low bias
- It can be difficult to interpret which terms in the model are most important

Ridge Regression

- When there are many correlated variables in a linear regression problem then the coefficients can be highly variable
 - The reason for this is that huge coefficients are balanced by equally small coefficients on correlated variables
- Solution: add a penalty for large coefficients

$$\hat{\beta}_{ridge} = \operatorname{argmin}_{\beta} \left\{ RSS + \lambda \sum_{j=1}^p \beta_j^2 \right\}$$

Ridge Regression

- Now the criterion for the solutions is to minimize $\Omega(\lambda)$:

$$\begin{aligned}\Omega(\lambda) &= (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) + \lambda\beta^T \beta \\ \frac{\partial \Omega}{\partial \beta} &= \frac{\partial \{(\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta)\}}{\partial \beta} + \frac{\partial (\lambda\beta^T \beta)}{\partial \beta} \\ 0 &= -2\mathbf{X}^T (\mathbf{y} - \mathbf{X}\beta) + 2\lambda\beta \\ (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})\beta &= \mathbf{X}^T \mathbf{y}\end{aligned}$$

$$\begin{aligned}\hat{\beta} &= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \\ \hat{\mathbf{y}} &= \mathbf{X}\hat{\beta} = \mathbf{X}(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}$$

Ridge Regression

- Another way to understand why this works is to look at the SVD of the solution:

Recall that the SVD of a matrix is:

$$\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T$$

with the properties that \mathbf{D} is a diagonal matrix and $\mathbf{U}^T\mathbf{U} = \mathbf{I}, \mathbf{V}\mathbf{V}^T = \mathbf{I}$

$$\begin{aligned}\hat{\mathbf{y}} &= \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}^T\mathbf{y} \\ \mathbf{X}\hat{\boldsymbol{\beta}} &= \mathbf{U}\mathbf{D}\mathbf{V}^T(\mathbf{V}\mathbf{D}^T\mathbf{U}^T\mathbf{U}\mathbf{D}\mathbf{V}^T + \lambda\mathbf{I})^{-1}\mathbf{V}\mathbf{D}^T\mathbf{U}^T\mathbf{y} \\ \mathbf{X}\hat{\boldsymbol{\beta}} &= \mathbf{U}\mathbf{D}\mathbf{V}^T(\mathbf{D}^2 + \lambda\mathbf{I})^{-1}\mathbf{V}\mathbf{D}^T\mathbf{U}^T\mathbf{y} \\ \mathbf{X}\hat{\boldsymbol{\beta}} &= \mathbf{U}\mathbf{D}(\mathbf{D}^2 + \lambda\mathbf{I})^{-1}\mathbf{V}^T\mathbf{V}\mathbf{D}^T\mathbf{U}^T\mathbf{y} \\ \mathbf{X}\hat{\boldsymbol{\beta}} &= \mathbf{U}\mathbf{D}(\mathbf{D}^2 + \lambda\mathbf{I})^{-1}\mathbf{D}\mathbf{U}^T\mathbf{y}\end{aligned}$$

$$\mathbf{X}\hat{\boldsymbol{\beta}} = \sum_{j=1}^p \mathbf{u}_j \frac{d_j^2}{d_j^2 + \lambda} \mathbf{u}_j^T \mathbf{y}$$

Ridge Regression

- Compare this to the least squares decomposition:

Ridge regression $\mathbf{X}\hat{\beta}_{ridge} = \sum_{j=1}^p \mathbf{u}_j \frac{d_j^2}{d_j^2 + \lambda} \mathbf{u}_j^T \mathbf{y}$

Least squares regression $\mathbf{X}\hat{\beta}_{LS} = \mathbf{U}\mathbf{U}^T \mathbf{y} = \sum_{j=1}^p \mathbf{u}_j \mathbf{u}_j^T \mathbf{y}$

- So we see that ridge regression shrinks the coefficients in the basis of \mathbf{U}
- Recall from PCA theory that this basis is related to components of maximum variance, and it appears λ shrinks components with small variance fastest
- Ridge regression is also known as \mathbb{L}_2 - regularization, Tikhonov regularization, or weight-decay in training neural nets

The Lasso

- Also known as \mathbb{L}_1 – regularization or basis pursuit
- The goal is the same as in ridge regression, but the penalty is different. Now we write the objective as

$$\hat{\beta}_{ridge} = \operatorname{argmin}_{\beta} \left\{ \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 \right\}, \text{ subject to } \sum_{j=1}^p |\beta_j| \leq t$$

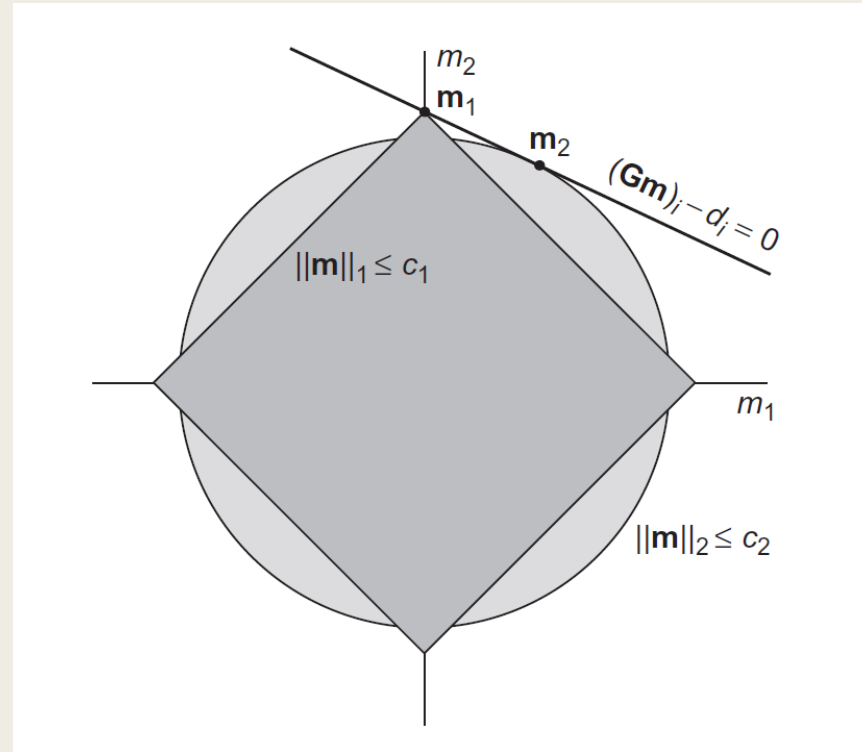
This can also be written in *Lagrangian form*:

$$\hat{\beta}_{ridge} = \operatorname{argmin}_{\beta} \left\{ \frac{1}{2} RSS + \lambda \sum_{j=1}^p |\beta_j| \right\}$$

The Lasso

- The Lasso is non-linear in \mathbf{y} so there is no closed-form solution. The optimal choice of β is found by iterative optimization procedures
- The advantage of the Lasso is that some coefficients can actually be set to 0 by the optimization procedure (as opposed to ridge regression, which only sends them to small values). This gives a clear interpretation of which terms are important.

Visualization of \mathbb{L}_2 vs \mathbb{L}_1 regularization



Regression using derived inputs

There are some methods of pre-processing our data which can achieve similar effects in reducing the model terms:

- PCA: we can find a regression along principle components instead of the original axes the data was given in.

Compute the SVD $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T$ and define $\mathbf{Z} = \mathbf{U}\tilde{\mathbf{D}}$ for $\tilde{\mathbf{D}} \in \mathbb{R}^{N \times M}$, $M < p$.

Now find the least squares solution β for $\mathbf{y} = \mathbf{Z}\beta^T$

- Partial least squares (PLS): similar to PCA in that we construct new components and perform the regression in a truncated version of this space.

Compute $\phi_{1j} = \langle \mathbf{x}_j, \mathbf{y} \rangle$ for each $j \in \{1, \dots, p\}$

Construct the derived input $\mathbf{z}_1 = \sum_j \phi_{1j} \mathbf{x}_j$

Compute the coefficient $\hat{\theta}_1 = \langle \mathbf{z}_1, \mathbf{y} \rangle / \langle \mathbf{z}_1, \mathbf{z}_1 \rangle$

Orthogonalize the remaining $\mathbf{x}_{j+1}, \mathbf{x}_{j+2}, \dots$ with respect to \mathbf{z}_1

Repeat the above procedure for all $\phi_{m,j}$ and \mathbf{x}_j that are left

Linear Methods of Classification

Linear Discriminant Analysis

Suppose we have a set of classes $G = \{1, 2, \dots, k, \dots, K\}$ and we wish to compute the probability that observation $X = x$ is in class $G = k$. This can be written using Bayes theorem as:

$$P(G = k|X = x) = \frac{P_k(x)P_G(k)}{\sum_{j=1}^K P_j(x)P_G(j)}$$

We assert that all classes are distributed as multivariate Gaussians of the form

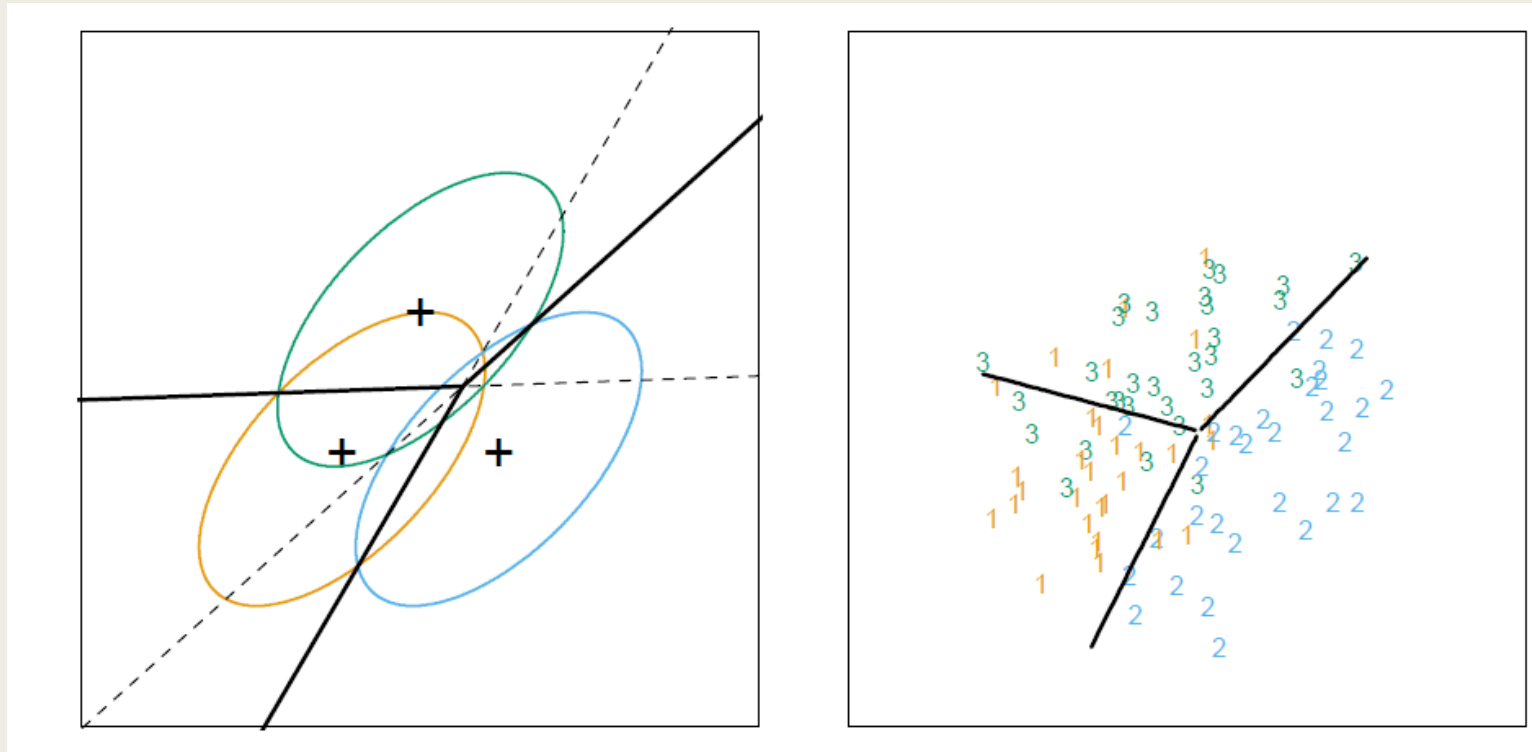
$$P_j(x) = \frac{1}{(2\pi)^{p/2} |\Sigma_k|^{1/2}} e^{-\frac{1}{2}(x-\mu_k)^T \Sigma_k^{-1} (x-\mu_k)}$$

If we assume Σ_k are equal for all k then we can compare the log-ratio of probabilities:

$$\log \frac{P(G = k|X = x)}{P(G = l|X = x)} = \log \frac{P_G(k)}{P_G(l)} + \log \frac{P_k(x)}{P_l(x)}$$

Linear Methods of Classification

Linear Discriminant Analysis



Logistic Regression

- We wish to model posterior log-probabilities of classes using linear functions in x , subject to the constraint that these functions sum to 1 and remain in $[0,1]$

$$\begin{aligned}\log \frac{\Pr(G = 1|X = x)}{\Pr(G = K|X = x)} &= \beta_{10} + \beta_1^T x \\ \log \frac{\Pr(G = 2|X = x)}{\Pr(G = K|X = x)} &= \beta_{20} + \beta_2^T x \\ &\vdots \\ \log \frac{\Pr(G = K - 1|X = x)}{\Pr(G = K|X = x)} &= \beta_{(K-1)0} + \beta_{K-1}^T x.\end{aligned}$$

Logistic Regression

- The consequence of this formulation is that

$$\begin{aligned}\Pr(G = k|X = x) &= \frac{\exp(\beta_{k0} + \beta_k^T x)}{1 + \sum_{\ell=1}^{K-1} \exp(\beta_{\ell 0} + \beta_{\ell}^T x)}, \\ \Pr(G = K|X = x) &= \frac{1}{1 + \sum_{\ell=1}^{K-1} \exp(\beta_{\ell 0} + \beta_{\ell}^T x)},\end{aligned}$$

Quick Derivation for 2 classes:

$$\log \frac{P(G = 1|X = x)}{P(G = 2|X = x)} = \beta_{10} + \beta_{11}x$$

$$e^{\beta_{10} + \beta_{11}x} = \frac{P(G = 1|X = x)}{P(G = 2|X = x)} = \frac{P(G = 1|X = x)}{1 - P(G = 1|X = x)}$$

$$P(G = 1|X = x) = \frac{e^{\beta_{10} + \beta_{11}x}}{1 + e^{\beta_{10} + \beta_{11}x}}$$

Logistic Regression

- Consider the two-class example. Use an indicator function to code the classes:

$$y_i = \begin{cases} 1 & \text{if } x_i \text{ is in class 1} \\ 0 & \text{if } x_i \text{ is in class 2} \end{cases}$$

We can write the log-likelihood of a set of N observations as a function of the parameters $\beta = \{\beta_{10}, \beta_1\}$

$$\begin{aligned} \ell(\beta) &= \sum_i \log p_k(x_i; \beta) = \sum_{i=1}^N \{y_i \log p(x_i; \beta) + (1 - y_i) \log(1 - p(x_i; \beta))\} \\ \ell(\beta) &= \sum_{i=1}^N \{y_i \beta^T x_i - \log(1 + e^{\beta^T x_i})\} \end{aligned}$$

Logistic Regression

- Find the maximum likelihood estimate of β by taking the derivative

$$\frac{\partial \ell(\beta)}{\partial \beta} = 0 = \sum_{i=1}^N x_i (y_i - p(x_i; \beta))$$

- Typically this set of p equations for β is solved using the Newton-Raphson method

$$\beta^{new} = \beta^{old} - \left(\frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta^T} \right)^{-1} \frac{\partial \ell(\beta)}{\partial \beta}$$

If we define \mathbf{W} a $N \times N$ diagonal matrix such that $w_{ii} = p(x_i; \beta^{old})(1 - p(x_i; \beta^{old}))$ then the Newton-Raphson step becomes

$$\beta^{new} = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} (\mathbf{X} \beta^{old} + \mathbf{W}^{-1} (\mathbf{y} - \mathbf{p}))$$

Note: this is also known as the iteratively reweighted least squares (IRLS) algorithm

That's all!

- For more details, I recommend *The Elements of Statistical Learning* by Hastie, Tibshirani, and Friedman