Things Linear (like algebra and models and stuff)

Schwartz

November 15, 2016

0 Synopsis

Glorified cheat sheet

1 LinAlg

1.1 inner product – $dot \ product$

 $||x||_2, ||x||_1, \text{ and } ||x||_{\infty}$

1.2 outer product (a.k.a. a whole lot of inner product)

$$\begin{bmatrix} \left| \right| \\ \left| \right| \\ \left| \right| \\ | \right| \\ | n \times 1 \end{bmatrix} = \begin{bmatrix} C_{r,c} = A_r \times B_c \\ \\ | \right| \\ | n \times m \end{bmatrix}$$

1.3 broadcasting

A little bit like outer product, but also a whole lot more...

$$\begin{bmatrix} \\ \\ \\ \\ \\ \end{bmatrix}_{n\times 1}^{A} * \begin{bmatrix} \\ \\ \\ \end{bmatrix}_{1\times m}^{B} = \begin{bmatrix} \\ \\ \\ \end{bmatrix}_{1\times m}^{B} * \begin{bmatrix} \\ \\ \\ \\ \end{bmatrix}_{n\times 1}^{A} = \begin{bmatrix} \\ \\ \\ \\ \end{bmatrix}_{n\times m}^{C} * B_{c} \end{bmatrix}_{n\times m}^{B}$$

1.4 matrix things and such 1.5

- MM is !communative but is associative and distributive
- $\bullet \ (AB)^T = B^T A^T$
- If $A = A^T$ then A is symmetric $\frac{1}{2}A + \frac{1}{2}A^T$ is symmetric
- trABC = trBCA = trCAB
- $\nabla_x a^T x = a$, and if A symmetric $\nabla_x x^T A x = 2Ax$ $\nabla_x^2 x^T A x = 2A$
- $(AB)^{-1} = B^{-1}A^{-1}$
- $(A^{-1})^T = (A^T)^{-1}$

1.6 Invertability

- column rank is row rank is rank
- $\operatorname{rank}(A_{(n,m)}) \leq \min(m,n), full \text{ if equal}$
- $A^{-1}A = I$ if $A \in \mathbb{R}^{n \times n}$ is full rank otherwise A is singular/non-invertable
- $(A^{-1})^T = (A^T)^{-1}$
- $U^T U = I (= U U^T)$ for orthonormal U
- $|\det A| = abs(|A|)$ is the volume of the \mathbb{R}^n -parallelotope formed by the vectors of A
- |A| = 0 if A is not full rank I.e., A singular (non-invertable)
- $|A^{-1}| = |A|^{-1}$, $|A| = |A^T|$, |AB| = |A||B|

• span $\{x_1, \dots, x_m\} = \mathcal{R} (A = [x_1, \dots, x_m])$ $= \left\{ v \in \mathbb{R}^n : v = \sum_{i=1}^m \alpha_i x_i \right\}$ with $x_i \in \mathbb{R}^n$ and $\alpha_i \in \mathbb{R}$

Relevant for LM's:

Spaces

The projection of y onto $\mathcal{R}(A = [x_1, \dots, x_m])$ is

$$\underset{v \in \text{span}\{x_1, \dots, x_m\}}{\operatorname{argmin}} ||y - v||_2 = A(A^T A)^{-1} A^T y$$

Relevant for SVM's:

The projection of y onto a is $\frac{a^T}{||a||}y = \frac{a^T}{\sqrt{a^Ta}}y$

- The nullspace of A ∈ R^{n×m} N(A) = {x ∈ R^m : Ax = 0}
 Why are R(A^T) and N(A) orthogonal complements?
 I.e., {w : w = u + v, u ∈ R(A^T), v ∈ N(A)} = R^m
 and R(A^T) ∩ N(A) = {0}
 - $A = A^T$ is positive semidefinite if $x^T A x \ge 0$ and positive definite if $x^T A x > 0$ \implies full rank (i.e., non-singular/nvertable)
 - For full rank $A_{(n,m)}$ with n > mthe *Gram matrix* A^TA is positive definite Relevant for LM's:

 X^TX is inverted in the least squares fit $X(X^TX)^{-1}X^Ty$

•
$$A^{-1} = \frac{1}{|A|} \operatorname{adj}(A) \in \mathbb{R}^{n \times n}$$

 $\operatorname{adj}(A)_{i,j} = (-1)^{i+j} |A_{-i,-j}|$

1.7 Eigens

• An eigenvector $(x \neq \mathbf{0} \text{ in } \mathbb{C}^n)$ and eigenvalue (λ) pair of $A \in \mathbb{R}^{n \times n}$ satisfy

$$Ax = \lambda x$$

• Solutions for $(\lambda I - A)x = 0$ exist if $(\lambda I - A)$ is singular/the nullspace $\mathcal{N}(\lambda I - A) \neq \{\mathbf{0}\}$

[so that rank
$$((\lambda I - A)^T)$$
 = rank $(\lambda I - A)$ is not full]

• In which case eigenvalue solutions to $|\lambda I - A| = 0$ can be used to solve for eigenvectors in

$$(\lambda I - A)x = 0$$

• And for eigenvalues $\lambda_1, \dots \lambda_n$ of A we have that

$$\operatorname{tr} A = \sum_{i=1}^{n} \lambda_{i}$$
 and $|A| = \prod_{i=1}^{n} \lambda_{i}$ and $\operatorname{rank}(A) = \sum_{i=1}^{n} 1_{[\lambda_{i} \neq 0]}$

- Quiz: what are the eigens for (i) diagonal matrix D and (ii) non-singular A?
- For $A = A^T$ we have that (i) $\lambda_1, \dots, \lambda_n \in \mathbb{R}^n$ and (ii) the eigenvectors are orthonormal

Relevant for MVN:

All the eigenvalues of Σ must be non zero so that Σ^{-1} exists

• For linearly independent eigenvectors $X = [x_1, \dots x_n]$ and associated eigenvalues $\Lambda = \operatorname{diag}(\lambda_1, \dots \lambda_n)$

$$AX = X\Lambda$$
 implies that $A = X\Lambda X^{-1}$ is diagonalizable

Relevant for PCA:

Diagonalization a.k.a. spectral or eigenvalue decomposition is a special case of the singular value decomposition which uses the covariance/correlation matrix rather than the data matrix

• Quiz: A^{-1} for diagonalizable A? What sign are the eigenvalues of positive definite A?

2 Regression

The range of the (full rank) n samples p predictors covariate (or design) matrix $\mathcal{R}(X)$ is the space (of dimension p < n) of all possible predictions of $y \in \mathbb{R}^n$

$$\mathcal{Y} = X\boldsymbol{\beta}$$
, for $\boldsymbol{\beta} \in \mathbb{R}^p$

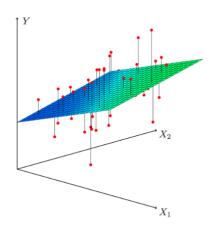
The projection of y onto $\mathcal{R}(X)$ is \hat{y} is

$$\min_{X\boldsymbol{\beta}} ||X\boldsymbol{\beta} - y||_2 = \min_{X\boldsymbol{\beta}} \left((X\boldsymbol{\beta} - y)^T (X\boldsymbol{\beta} - y) \right)$$
$$= \min_{X\boldsymbol{\beta}} \left(\boldsymbol{\beta}^T X^T X \boldsymbol{\beta} - 2X^T y \boldsymbol{\beta} + y^T y \right)$$

which (by taking the derivative) is maximized at

$$\boldsymbol{\beta} = \left(X^T X\right)^{-1} X^T y$$

2.1 Linear Models



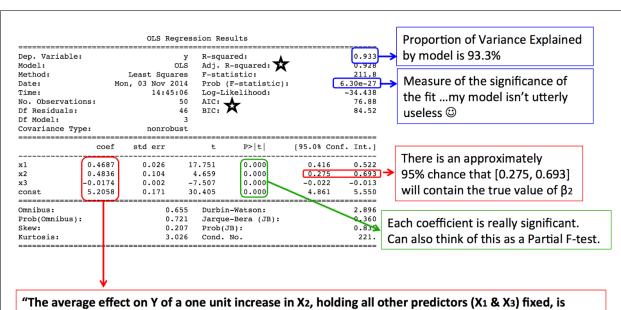
$$y_i = \beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip} + \epsilon_i, \ \epsilon_i \sim N(0, \sigma^2)$$

 $\hat{y}_i = \beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip}$

- Assumptions when do they matter?
- Higher order terms what does "linear" mean?
- Challenges interpreting coefficients?
- The role of rank(X) and multicollinearity?
- Significance testing and model building?

2.2 Fit

Total Variation	Total Error	Average Error	Proportion Modeled
$TSS = \sum_{i=1}^{n} (y_i - \bar{y})^2$	$RSS = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$	$RSE = \sqrt{\frac{1}{n-p-1}RSS}$ $= \sqrt{\frac{\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}{n-p-1}}$	$R^2 = \frac{TSS - RSS}{TSS}$ $= 1 - \frac{RSS}{TSS}$



"The average effect on Y of a one unit increase in X_2 , holding all other predictors ($X_1 \& X_3$) fixed, is 0.4836"

- However, interpretations are generally pretty hazardous due to correlations among predictors.
- p-values for each coefficient ≈ 0, so might be okay here

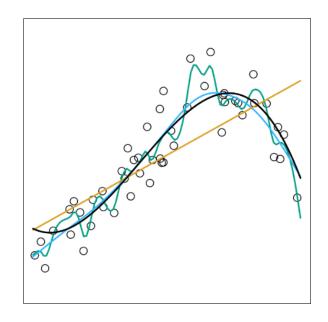
Note: Magnitude of the Beta coefficients is NOT how to determine whether predictor contributes. Why?

3 Overfitting

Let $y_0 = \theta + \epsilon_0$ with $\theta = f(x_0)$ and $\epsilon \sim N(0, \sigma_{\epsilon}^2)$ For estimator $\hat{\theta} = \hat{f}(x_0)$,

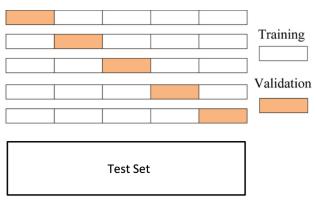
$$MSE = \frac{1}{n} \sum_{i} (y_i - \hat{\theta})^2 \approx E \left[(y_i - \hat{\theta})^2 \right]$$
$$= \sigma_{\hat{\theta}}^2 + (E[\hat{\theta}] - \theta)^2 + \sigma_{\epsilon}^2$$
$$= Variance + Bias^2 + Noise$$

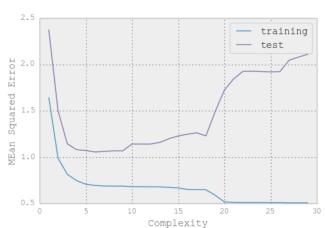
$$MSE(\hat{\theta}) = E\left[\left(\hat{\theta} - \theta\right)^{2}\right]$$
$$= \sigma_{\hat{\theta}}^{2} + \left(E[\hat{\theta}] - \theta\right)^{2}$$
$$= Variance + Bias^{2}$$



The variance/bias tradeoff

3.1 K-folds cross-validation





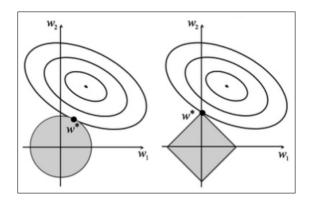
3.2 Regularization

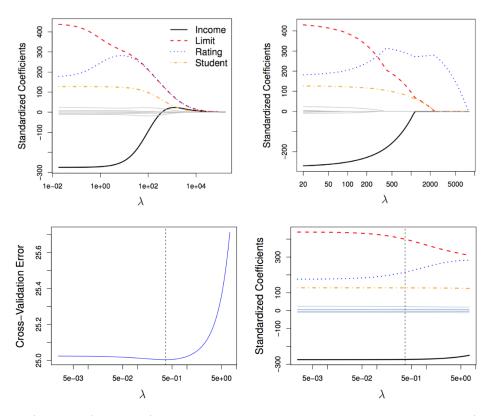
$$LS = \min_{\beta} ||X\beta - y||_2 \qquad (no penalty)$$

Ridge =
$$\min_{\beta} ||X\beta - y||_2 + \lambda ||\beta||_2$$
 (L₂penalty)

Lasso =
$$\min_{\beta} ||X\beta - y||_2 + \lambda ||\beta||_1$$
 (L1penalty)

Penalized coefficients not scale equivariant!



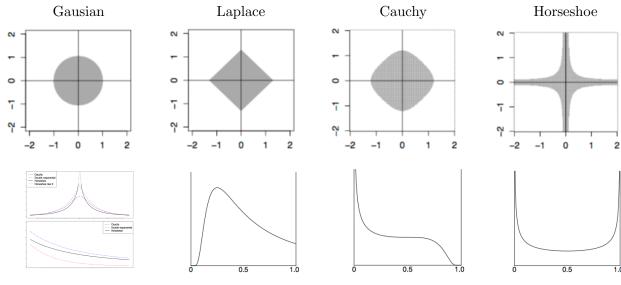


Ridge (top left), Lasso (top right), and cross-validation tuning parameter selection (bottom row)

3.3 Summary

1. Model selection 2. Regularization 3. Dimension reduction (e.g., principal components reduction)

3.4 Bonus: Bayesian regularization priors



Tails Implied shrinkage prior profiles from none (0) to total (1) shrinkage

4 Logistic Regression

- The logit link function $g(p) = \log\left(\frac{p}{1-p}\right)$ maps $p \in [0,1] \mapsto Z \in \mathbb{R}$
- For a binary outcome Y, setting $E[Y] = g^{-1}(Z) = \frac{\exp(Z)}{1 + \exp(Z)}$ and $Z = \beta_0 + \beta_1 X_1 + \dots + \beta_m X_m$

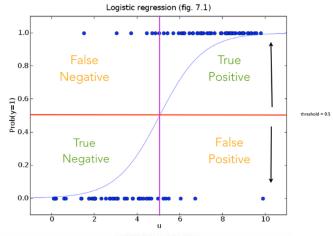
$$Pr(Y = 1|X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \dots + \beta_n X_m)}}$$

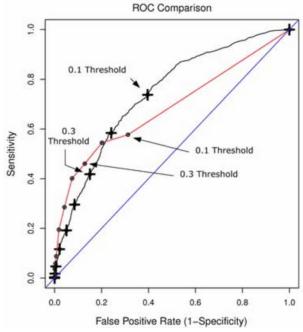
and

$$\exp(\beta_0)\exp(\beta_1 X_1)\cdots\exp(\beta_m X_m) = \frac{Pr(Y=1|X)}{Pr(Y=0|X)}$$

with $\exp(\beta_j)$ the multiplicative (logarithmic scale) odds increase for a 1-unit increase in X_j

- True Positive (TP)
- True Negative (TN)
- False Positive (FP) (Type I error)
- False Negative (FN)(Type II error, 1 - power)
- Accuracy: (TP + TN)/(TP + FP + FN + TN)
- F1 score: 2TP/(2TP + FP + FN)
- Sensitivity: TP/(TP + FN) (power)
- Specificity: TN/(TN + FP)(1 - Type I error rate α)
- Positive Predictive Value: TP/(TP + FP) (Precision)
- Negative Predictive Value: TN/(TN + FN)
- False Positive Rate: FP/(FP + TN) (fall-out)
- False Discovery Rate: FP/(FP + TP)(1 - precision)
- False Negative Rate: FN/(FN + TP)(1 - sensitivity)





Appendix: classical model selection

$$C_p = \frac{1}{n}(RSS + 2p\hat{\sigma}^2) \longleftarrow \begin{array}{l} \text{Mallow's Cp} \\ \text{p is the total \# of parameters} \\ \hat{\sigma}^2 \text{ is an estimate of the variance of the error, } \epsilon \end{array}$$

$$AIC = -2logL + 2 \cdot \underline{p} \ \ \, \ \ \, \ \ \,$$
 L is the maximized value of the likelihood function for the model estimated

$$BIC = \frac{1}{n}(RSS + log(n)\underline{p}\hat{\sigma}^2) \longleftarrow \text{This is AIC, except 2 is replaced by log(n).} \\ \log(n) > 2 \text{ for n>7, so BIC generally exacts a heavier penalty for more variables}$$

Can show AIC and Mallow's Cp are equivalent for linear case