

Model Assumptions

Linearity Assumption - different than GLM/MLR

Linear relationship between predicting variables x_1, \dots, x_p and $\log(\text{E}(Y))$ aka $\log(\text{lognormal})$

- Plot: Predictions vs Residuals
- Plot: Predictions vs log(|y|) of outcome data

Independence Assumption - different than GLM/MLR

Y_1, \dots, Y_n are independent RV

- Plot: Predictions vs Residuals - cannot check independence, settle for uncorrelated residuals

Normality Assumption

Normally distributed residuals $\epsilon_1, \dots, \epsilon_n$ aka $\log(\text{normal})$

Normality Assumption - CANONICAL REGRESSION IS IDIOSYNCRATIC! need to check normality of residuals (Pearson residuals or deviance residuals - should follow approximately N(0,1) if the model is a good fit)

- Plot: Q-Q plot
- Plot: Histogram

$Y_j | (x_{j1}, \dots, x_{jp}) \sim \text{Poisson}(\lambda_{j1}, \dots, \lambda_{jp})$

Estimated rates $\lambda_i = \lambda_i(x_{i1}, \dots, x_{ip}) = e^{\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}}$

Distributions - see Poisson section

For test on ALL regression coefficients (not β_0 if NOT an F test, it's Deviance test)

Full Model is $\log(\text{E}(Y)) = \beta_0 + \beta_1 x_{11} + \dots + \beta_p x_{1p} + \beta_{p+1} + \beta_{p+2} + \dots + \beta_{p+q}$

Reduced Model is $\log(\text{E}(Y)) = \beta_0 + \beta_1 x_{11} + \dots + \beta_p x_{1p}$

H_0 : at least one $\beta_j = 0$

Produce χ^2_{p+q-1} Deviance test, reject H_0 if $p < \alpha$

- null deviance - Residual deviance deviance_testing deviance_res ~ null log likelihood - full log likelihood
- null deviance $\sim \chi^2_{p+q-1}$ Residual deviance deviance_testing deviance_res $\sim \chi^2_{p+q-1}$ deviance_testing deviance_res
- Compute p-value for Deviance test statistic: $1 - \text{pnorm}(\text{deviance_testing deviance_res}, df = p + \text{predictors})$

For test on SUBSET of regression coefficients (it's NOT an F test, it's Deviance test)

Full Model is $\log(\text{E}(Y)) = \beta_0 + \beta_1 x_{11} + \dots + \beta_p x_{1p} + \beta_{p+1} + \beta_{p+2} + \dots + \beta_{p+q}$

Reduced Model is $\log(\text{E}(Y)) = \beta_0 + \beta_1 x_{11} + \dots + \beta_p x_{1p}$

H_0 : at least one $\beta_j = 0$

Produce χ^2_{p+q-1} Deviance test, reject H_0 if $p < \alpha$

- Deviance of reduced model - Deviance of the full model: deviance_testing deviance_res - reduced log likelihood - full log likelihood
- Deviance test statistic: χ^2_{p+q-1}
- Compute p-value for Deviance test statistic: $1 - \text{pnorm}(\text{deviance_testing deviance_res}, df = p + \text{predictors in subset})$

Somehow related: Goodness of fit Hypothesis Testing (aka Pearson or Deviance residuals for this Deviance test)

H_0 : model does NOT fit data

Deviance test and sum of deviance residuals $\sim \chi^2_{p+q-1}$

OR Deviance test and sum of deviance residuals $\sim \chi^2_{p+q-1}$

Produce χ^2_{p+q-1} Deviance test, reject H_0 if $p < \alpha$

THIS IS MORE WHEN WANT LARGE P-VALUES (DO THAT WE DO NOT REJECT H_0 MEANING MODELS FITS DATA)

For test of statistical significance on one predictor, given all other predictors in model: Wald test

See section Hypothesis testing using a test statistic for all 3 (left-tailed, 2-tailed, right-tailed)

$H_0: \beta_j = 0$

Test stat: $\frac{\hat{\beta}_j}{\text{SE}(\hat{\beta}_j)}$ reject H_0 if $p < \alpha$

Overdispersion in GLM - see Dr Serban notes

Correcting Training Risk for bias in Linear Models:

Given R is number of predictors in the model and $R(S) = R_{\text{train}}(S) + \text{Some Complexity Penalty CP}$

- using Mallows's CP: $R(S) = R_{\text{train}}(S) + \frac{2 \cdot \text{MSE}_{\text{train}}(S)}{n}$, where $\text{MSE}_{\text{train}}(S)$ is estimated variance of full model (not always possible to estimate e.g. when $p = n$)
- using Akaike's IC CP: $R(S) = R_{\text{train}}(S) + \frac{2 \cdot \text{MSE}_{\text{train}}(S)}{n}$, where $\text{MSE}_{\text{train}}(S)$ is true variance of full model, not estimated
- need to replace $\text{MSE}_{\text{train}}(S)$ with $\text{MSE}_{\text{train}}(S) \approx \text{MSE}_{\text{train}}(S)$
- Important! Most software routines $\text{MSE}_{\text{train}}(S) \approx \text{MSE}_{\text{train}}(S)$
- Akaike Information Criterion is an estimate for the predictor risk
- using Bayesian IC CP: $R(S) = R_{\text{train}}(S) + \frac{\text{MSE}_{\text{train}}(S)}{n}$
- need to replace $\text{MSE}_{\text{train}}(S)$ with $\text{MSE}_{\text{train}}(S) \approx \text{MSE}_{\text{train}}(S)$
- BIC: produces completely more than other approaches \Rightarrow preferred in model selection
- Leave-One-Out CV Approximation: $R(S) = R_{\text{train}}(S) + \frac{2 \cdot \text{MSE}_{\text{train}}(S)}{n}$
- Because $R_{\text{train}}(S) \leq \text{MSE}_{\text{train}}(S) \leq \text{MSE}_{\text{train}}(S)$ LOO CV Approximation CP penalizes completely less than Mallows's CP

Leave-One-Out CV (if direct measurement of predictive power): $R(S) = \frac{1}{n} \sum_{i=1}^n (\hat{y}_{-i} - y_i)^2$ where \hat{y}_{-i} is \hat{y} estimated from a model fitted without observation i LOO CV is approximately AIC when n is replaced by $\text{MSE}_{\text{train}}(S)$

Correcting Training Risk for bias in Generalized Linear Models (e.g. in Logistic or Poisson Regression):

Training risk $R_{\text{train}}(S)$ for a submodel S , third argument for submodel S is future observation Y

$R_{\text{train}}(S) = \frac{1}{n} \sum_{i=1}^n \log \left(\frac{e^{\hat{\eta}_i}}{1 + e^{\hat{\eta}_i}} \right) + \frac{1}{n} \sum_{i=1}^n \log \left(\frac{e^{\hat{\eta}_i}}{1 + e^{\hat{\eta}_i}} \right)$ sum of square deviances of submodel S

- Akaike Information Criterion CP and Bayesian Information Criterion CP are commonly used for model selection in GLM since they are defined in terms of log likelihood function

Penalizing the [Minimum] Sum of Squared Errors (aka Sum of Least Squares) in Regularized Regression

Inside the $\hat{\beta}$ in my physical notes are meant to be $\hat{\beta}_j$ and $\hat{\beta}_j$ are in $\hat{\beta}$

$\text{RSS}(\hat{\beta}) = \sum_{i=1}^n (\hat{y}_i - y_i)^2 = \sum_{i=1}^n (\hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_p x_{ip})^2 + \sum_{i=1}^n (\hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_p x_{ip})^2 + \sum_{i=1}^n (\hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_p x_{ip})^2$

The penalty constant λ in $\hat{\beta}$ penalized or regularized regression controls the trade-off between lack of fit and model complexity

Choices of Penalty:

- L_1 penalty: $\|\hat{\beta}\|_1 = \sum_{j=1}^p |\hat{\beta}_j|$ $\hat{\beta}_j = 0$! produces best model given a selection criterion, but it requires fitting all possible submodels which may not be possible

- L_2 penalty: Lasso $\|\hat{\beta}\|_2 = \sqrt{\sum_{j=1}^p \hat{\beta}_j^2}$ minimizes sparsity

- L_1 penalty, Ridge $\|\hat{\beta}\|_2 = \sqrt{\sum_{j=1}^p \hat{\beta}_j^2}$ easy to implement but it does not do variable selection, does not distinguish between sparse and non-sparse vectors

Note! For performing model selection using regularized regression, predicting variables MUST be **centered** (standardized, response is recommended to be **uncentered**). After selecting the "best" model, use the original scale when fitting the selected model for interpretation of the regression coefficients.

Note! The L_1 penalty produces sparse estimates, while the L_2 penalty does not.

$\sum_{j=1}^p \hat{\beta}_j = 0$, $\sum_{j=1}^p \hat{\beta}_j^2 = 1$

$\sum_{j=1}^p \hat{\beta}_j = 0$, $\sum_{j=1}^p \hat{\beta}_j^2 = 1$

In this class

SLR/MLR

goodness of fit assumption means analyzing residuals from SLR or regular residuals, for MLR's constant variance - standardized (the lecture slides forget to do that, for MLR's other assumptions - regular) to check if the 4 model assumptions hold.

It requires using **goodness of fit** (aka **goodness of fit**) through cross-validation will be more on predictive power.

the overall F test, the partial F test, and the individual p-values of the coefficients to determine statistical significance of all predictors including intercept, ϵ is a subset of predictors, and of an individual predictor, respectively.

Logistic Regression

goodness of fit assumption means

- analyzing residuals (just deviance) to check if the 3 model assumptions hold (linearity, independence - more the uncorrelated errors, the function, AIC checking those residuals are normally distributed (even though it's not in the original assumptions))

- AND performing hypothesis testing on residuals (both Pearson and deviance), i.e. computing the p-value of the Deviance statistic D (where a LARGE p-value is NOT REJECT H_0 , meaning model fits data well).

predictive power means comparison of classification error from cross-validation at different thresholds

Poisson Regression

goodness of fit assumption means

- analyzing residuals (just deviance) to check if the model assumptions hold (linearity, independence - more the uncorrelated errors, 77 variance assumption), AND checking those residuals are normally distributed (even though it's not in the original assumptions)

- AND performing hypothesis testing on residuals (both Pearson and deviance), i.e. computing the p-value of the Deviance statistic D (where a LARGE p-value is NOT REJECT H_0 , meaning model fits data well).

predictive power means comparison of classification error from cross-validation at different thresholds

From Data Split

Question on testing subsets of variables and model comparison:

- Residual ratio tests compare models
- Tests the form of partial F test in multiple linear regression and simple linear regression

- Difference of the log of the differences in any generalized linear model
- can compare any models with this concept

Model Selection

We cannot perform variable selection based on the statistical significance of the regression coefficients, statistical significance is only true in the context of the given model. Also, a better model can be found, even if the current one is statistically significant. It's possible to select a model to include variables that are not statistically significant, even though that model still provides the best prediction, for example, and vice versa.

Once for model selection (split a lot of variables, use model using grid (or knn) with the variables on their original scale

- Best Subset (not possible if p is large, 2^p subsets to check)
- Greedy algorithm (forward stepwise, backward stepwise, forward backward stepwise)

- Backward and forward stepwise regression will generally provide different sets of selected variables when p , the number of predicting variables, is large. Backward stepwise cannot be performed if p is larger than n

- Forward stepwise regression is preferable over backward stepwise regression because it starts with smaller models.

- Global algorithm - regularized regression (lasso, ridge (not for selection), elastic net) - lasso - ridge
- Better fitting. Most standardizing predicting variables (standardized error), recommended to standardize response

Note on scaling vs standardization

Standardization is transforming your data so it has mean 0 and standard deviation 1, like a standard normal distribution: $z = (x - \mu) / \sigma$

Scaling is transforming your data so it's 0-1 range: $x = (x - \min(x)) / (\max(x) - \min(x))$

Scaling means dividing variable by its standard deviation.

Combination of the two is called standardization.

	Lasso	Ridge
Does not work when (number of predictors p is larger than n)	Lasso performs variable selection (selects a subset of predictors)	Ridge does not perform variable selection
There is NO closed form regression for estimated regression coefficients, generalized additive models	Lasso does not deal with multicollinearity (a subset of high collinear predictors will be dropped, zero, resulting in $\hat{\beta}$ coefficients)	Ridge deals with multicollinearity (a subset of high collinear predictors will be dropped, zero, resulting in $\hat{\beta}$ coefficients)
Estimated regression coefficients are less biased than those from OLS (even model selection and bias, use OLS to estimate coefficients)	Lasso has some coefficients equal to 0	Ridge has some coefficients equal to 0

Elastic Net = Lasso + Ridge

$\sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^p \hat{\beta}_j^2 + \lambda_1 \sum_{j=1}^p |\hat{\beta}_j|$ where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, $\lambda_1 = 0$ lasso, $\lambda_1 = \infty$ ridge

λ_1 is the L_1 norm, λ_2 is the L_2 norm, where $\lambda_1 > 0$ elastic net, λ