

gamdist: Generalized Additive Models in Python

Bob Wilson

September 30, 2018

Abstract

TBD

Contents

1	Introduction	3
2	Generalized Additive Models	3
3	The Alternating Direction Method of Multipliers	6
4	Software Architecture	6
A	Properties of the Linear Model	6
A.1	Properties of the Estimated Model	6
A.2	Confidence Intervals and Hypothesis Tests	7
A.3	Model Selection	10
A.4	Checking Model Assumptions	12

1 Introduction

This paper introduces a Python library called `gamdist`, which uses a distributed optimization technique called the Alternating Direction Method of Multipliers (ADMM) to fit a special type of regression model called a Generalized Additive Model (GAM) to data.

Outline of Paper In §2 we describe Generalized Additive Models. In §3 we describe the Alternating Direction Method of Multipliers and how it may be used to fit GAMs. In §4, we describe the architecture of the library, including relevant implementation details.

2 Generalized Additive Models

The primary goal of `gamdist` is the estimation of certain aspects of the joint distribution of a collection of one or more random variables X called *features* and a random variable Y we will call the *response*. Specifically, we are interested in the conditional distribution of $Y | X$. Perhaps the simplest approach is the linear model, which assumes $Y | X = x \sim \mathcal{N}(\mu(x), \sigma^2)$, where $\mu(x) = \nu(x)^T \beta$. This notation captures three key assumptions of the linear model. First, the conditional distribution is Gaussian for all values of X . Second, the mean of the distribution depends on the features X in a fairly specific way discussed below. Third, the variance is the same for all values of X . These assumptions are all loosened in various generalizations of the linear model used in `gamdist`.

If we choose $\nu(x) = x$, then the assumption is that $\mu(x) = x^T \beta$, and the mean depends linearly on the features; however, the *linear* in *linear model* refers to the dependence on β , not on the features. It is common to include a constant term in $\nu(x)$ to account for an affine dependency between the features and the response. For example, $\nu(x) = [1 \ x_1 \ x_2 \ \dots]^T$. We might incorporate quadratic terms to capture nonlinear dependencies including interactions, such as $\nu(x) = [x_1 \ x_2 \ x_1^2 \ x_1 \cdot x_2 \ x_2^2 \ \dots]^T$. Or we might include more exotic transformations of the features, like $\nu(x) = [\log(x_1) \ \sin(x_2) \ \dots]^T$. These models are non-linear in the features, but linear in the parameters β ; however, the linear model will only incorporate transformations that we explicitly include.

Another common situation is when some or all of the features are categorical. For example, consider a model with a single feature corresponding to a person’s favorite color, and suppose choices are limited to red, green, and blue. The model would consist of the average responses for people who prefer any particular color. We might support such a model by defining

$$\nu(x) = \begin{cases} \begin{bmatrix} 1 & 0 & 0 \end{bmatrix}^T & \text{if } x = \text{red} \\ \begin{bmatrix} 0 & 1 & 0 \end{bmatrix}^T & \text{if } x = \text{green} \\ \begin{bmatrix} 0 & 0 & 1 \end{bmatrix}^T & \text{if } x = \text{blue}. \end{cases}$$

Alternative approaches to encoding categorical variables are common and useful in different

circumstances. We see that the simple linear model is applicable to a wide range of problems, even those that may not appear linear at first glance.

Fitting a linear model to a set of observations is called linear regression, and is accomplished by solving a least squares optimization problem. This is an example of maximum likelihood estimation (MLE), itself a special case of maximum a posteriori (MAP) estimation, which is the unifying approach used throughout `gamdist`. It is worth formulating this optimization problem so that we may see how it evolves as we consider more general scenarios.

Suppose we have n observations of the form $\{(x^{(i)}, y^{(i)})\}_{i=1,\dots,n}$. We assume these observations are independent and drawn from the same (unknown) joint distribution.¹ Under the assumptions of the linear model, $Y \mid X = x^{(i)} \sim \mathcal{N}(\mu(x^{(i)}), \sigma^2)$. The likelihood of a particular observation is

$$\mathcal{L}(\beta; x^{(i)}, y^{(i)}) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp\left(-\frac{1}{2\sigma^2} \left(y^{(i)} - \nu(x^{(i)})^T \beta\right)^2\right).$$

Note that by convention the likelihood is interpreted as a function of β parameterized by the observation $(x^{(i)}, y^{(i)})$. The likelihood of the entire set of observations is the product of the likelihoods of the individual observations: $\mathcal{L}(\beta; x, y) = \prod_{i=1}^n \mathcal{L}(\beta; x^{(i)}, y^{(i)})$. The log-likelihood is the sum of the log-likelihoods of the individual observations:

$$\ell(\beta; x, y) = \log \mathcal{L}(\beta; x, y) = -n/2 \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n \left(y^{(i)} - \nu(x^{(i)})^T \beta\right)^2.$$

Maximizing the likelihood is the same as maximizing the log-likelihood, and if we are only interested in estimating β , this is equivalent to the problem

$$\text{minimize} \quad \sum_{i=1}^n \left(y^{(i)} - \nu(x^{(i)})^T \beta\right)^2,$$

where the variable is β and $y^{(i)}$ and $\nu(x^{(i)})$ are data. An elementary result in optimization theory is that a unique solution exists if and only if $V^T V$ is full rank, where the i th row of V is equal to $\nu(x^{(i)})^T$. In that case, the optimal β satisfies the so-called normal equations:

$$V^T V \hat{\beta} = V^T y.$$

If we assume the model is correct (that is, that the conditional distribution really has the assumed form), exact formulae exist for confidence intervals on the parameters β . If the variance is unknown, it too can be estimated from the data. We can employ hypothesis tests against the null hypothesis that some or all of the components of β are zero. We can apply the resulting model to new data assumed to be drawn from the same joint distribution to compute confidence intervals on the response, $Y \mid X = x_{\text{new}}$ or the mean of this distribution,

¹This assumption is loosened in random effects or mixed effects models which are not considered here; see [Str12].

$\mu(x_{\text{new}})$. These are immensely valuable tools in the analysis of data and the application of data to predictions. A good reference on linear models is [Wei05]. Useful results are collected in Appendix A.

We can loosen the assumption of constant variance slightly without materially affecting the inferences that may be drawn. Suppose the i th observation has variance $\sigma^2/\phi^{(i)}$, where $\phi^{(i)}$ is a known quantity, but σ^2 may or may not be known a priori.

Even if the assumptions underlying the linear model are correct, if the noise is high, or if the number of features is large relative to the number of observations, we may use regularization to improve both the estimates of β and predictions based on the estimated model. Regularization reduces the sensitivity of the estimates to noise at the expense of introducing bias, and may be thought of as imposing a Bayesian prior on the parameters β . Some of the most common forms of regularization include ridge regression and the lasso [Tib96]. For example, the lasso may be formulated as the problem:

$$\text{minimize} \quad \sum_{i=1}^n \left(y^{(i)} - \nu(x^{(i)})^T \beta \right)^2 + \lambda \cdot \|\beta\|_1,$$

but this problem does not have a closed-form solution. Moreover, introducing regularization means the formulation is no longer a maximum likelihood estimation problem. Instead, it is a maximum a priori estimation problem. MLE problems have some statistical properties that MAP estimation problems do not possess.

If the constant variance assumption does not hold, but the dependence is known, then Weighted Least Squares.

If the conditional distribution is not Gaussian, other techniques may prove more useful. For example, if the conditional distribution is Laplacian, we may use least absolute deviation regression instead of least squares [BD93]. Like with the lasso, there are no exact formulae for statistical inference in this context and we must settle for an asymptotic or non-parametric approach such as the bootstrap [ET93].

Yet another approach to extending linear models was introduced by [NW72] and discussed in detail in [MN89]. Their formulation extends the linear model in a few ways. The conditional distribution is not assumed to be Gaussian. Common alternatives include the binomial and Poisson distributions. The mean of the distribution is permitted to depend on the features in a more complicated way, via the introduction of a *link function*, g : $\mu(x) = g^{-1}(\eta(x))$, where $\eta(x) = \nu(x)^T \beta$. When $g(x) = x$, this recovers the same relationship between the features and μ assumed in the linear model, but other link functions may be used like the logistic function $g(x) = \log(x/(1-x))$. Finally, the variance is sometimes permitted to depend on x instead of being constant. Such models are called Generalized Linear Models (GLMs). Fitting such models is accomplished via maximum likelihood estimation:

$$\text{minimize} \quad \sum_{i=1}^n \ell(\beta; x^{(i)}, y^{(i)}) + r(\beta),$$

where $r(\beta)$ is a regularization term on β , such as in the lasso. For many choices of distribution family and link function, the corresponding optimization problem is convex.

Introduced by [HT86], Generalized Additive Models (GAMs) extend GLMs by permitting $\eta(x)$ to be a nonparametric function of the features: $\eta(x) = \sum_{i=1}^p h_i(x_i)$, where h_i are smooth functions. When $h_i(x_i) = \beta_i x_i$, the linear model is recovered (all of the parametric dependencies discussed with regards to ν are still possible here of course, but the idea is that the data itself should tell us the form of the relationship). Typically the h_i functions are chosen to be some sort of spline, such as a natural cubic spline. GAMs are also fit via MLE, and many practical problems can be formulated as convex optimization problems.

(Reference GAM, GAMr, Casella and Berger, Seber) regularized GAMs

3 The Alternating Direction Method of Multipliers

4 Software Architecture

Acknowledgments

This is an example of an unnumbered section.

A Properties of the Linear Model

In this section, we state various useful facts (without proof) regarding the linear model. For details, see [Wei05], [SL03], [Woo17], and [CB01]. The goal of this section is twofold: to gather formulae used in `gamdist`, and to highlight what we typically want to do in a regression analysis, beyond just fitting the model. In fact, it is often desirable to quantify the uncertainty in fitted model parameters, check the assumptions of the model through examination of the residuals, perform model selection, and quantify the uncertainty associated with predictions.

The beauty of the linear model is that exact formulae are attainable in support of these goals. Other types of regression only support asymptotic or approximate formulae. However, it is rarely the case that the assumptions of the linear model are expected to hold exactly. If the assumptions are approximately valid, we may hope the formulae which follow are approximately valid. Much has been written about the robustness of these formulae under deviations from the assumptions [SL03, §9]. We are inspired by the maxim [BD87], “All models are wrong, but some are useful.”

A.1 Properties of the Estimated Model

Suppose $Y \mid X = x \sim \mathcal{N}(\mu(x), \sigma^2)$, where $\mu(x) = \nu(x)^T \beta$ and $\nu(x) \in \mathbf{R}^p$ is a known function. Let $\{(x^{(i)}, y^{(i)})\}_{i=1, \dots, n}$ be IID samples drawn from the joint distribution of X and Y . Let V be the matrix whose i th row is $\nu(x^{(i)})$, and assume V is full rank. Let $\hat{\beta} = (V^T V)^{-1} V^T y$. Then $\hat{\beta} \sim \mathcal{N}(\beta, \mathcal{I}^{-1})$, where $\mathcal{I} = (V^T V)/\sigma^2$ is the Fisher information matrix. Specifically, $\hat{\beta}$ has a multivariate normal distribution, and $\hat{\beta}$ is an unbiased estimate of β . It is also a

consistent estimate of β , meaning that $\hat{\beta}$ converges in probability to β , as n increases without limit. It is also the best linear unbiased estimate of β : any other linear, unbiased estimates have higher variance than $\hat{\beta}$ [Woo17, § 1.3.9].

A simple modification permits the model to be much more flexible. Suppose that the i th observation has variance $\sigma^2/w^{(i)}$, where $w^{(i)}$ are known, positive numbers. Let U be the matrix whose i th row is $\sqrt{w^{(i)}}\nu(x^{(i)})$, and let $z^{(i)} = \sqrt{w^{(i)}}y^{(i)}$. Then the conclusions of this section are valid substituting $y \rightarrow z$ and $V \rightarrow U$. For example, $\hat{\beta} = (U^T U)^{-1} U^T z$ is the best linear unbiased estimate of β [Wei05, § 5.1]. When $w^{(i)} = 1$, we get the original results since $V = U$ and $y = z$. In what follows, we will proceed in terms of V and y .

If σ^2 is unknown, it may be estimated from the data. Let

$$\hat{\sigma}^2 = \frac{\|y - V\hat{\beta}\|_2^2}{n - p}. \quad (1)$$

Then $(n - p)\hat{\sigma}^2/\sigma^2 \sim \chi_{n-p}^2$, and $\hat{\beta}$ and $\hat{\sigma}^2$ are independent [Wei05, § 3.4.4]. This indicates that $\hat{\sigma}^2$ is an unbiased estimate of σ^2 .

A.2 Confidence Intervals and Hypothesis Tests

Suppose $\mathbf{Prob}\{\chi_{n-p}^2 \in (\ell, u)\} = \alpha$; that is, (ℓ, u) is a confidence interval (at level α) on a χ_{n-p}^2 random variable. Then

$$\mathbf{Prob}\left\{\sigma^2 \in \left(\frac{(n - p) \cdot \hat{\sigma}^2}{u}, \frac{(n - p) \cdot \hat{\sigma}^2}{\ell}\right)\right\} = \alpha. \quad (2)$$

A particularly useful case is when $u \rightarrow \infty$, corresponding to $\ell = \Phi_{\chi_{n-p}^2}^{-1}(1 - \alpha)$, where $\Phi_{\chi_{n-p}^2}$ is the cumulative distribution function of a χ_{n-p}^2 random variable, in which case

$$\mathbf{Prob}\left\{\sigma^2 < \frac{(n - p) \cdot \hat{\sigma}^2}{\Phi_{\chi_{n-p}^2}^{-1}(1 - \alpha)}\right\} = \alpha,$$

corresponding to an upper confidence limit on σ^2 , at level α .

We may compute confidence intervals on linear combinations of the components of β by noting that $c^T \hat{\beta} \sim \mathcal{N}(c^T \beta, \|c\|_{\mathcal{I}}^2)$, where $\|c\|_{\mathcal{I}}^2 = c^T \mathcal{I}^{-1} c$ is (the square of) the Mahalanobis norm [SL03, § 3.11.1], and thus

$$\frac{c^T \hat{\beta} - c^T \beta}{\|c\|_{\mathcal{I}}} \sim \mathcal{N}(0, 1). \quad (3)$$

Let $z^{1-\alpha/2}$ be the upper $\alpha/2$ quantile of a standard Gaussian random variable. Then

$$c^T \hat{\beta} \pm z^{1-\alpha/2} \cdot \|c\|_{\mathcal{I}} \quad (4)$$

are the endpoints of a $100(1 - \alpha)\%$ confidence interval on $c^T \beta$. This formula is only computable when \mathcal{I} is computable, which in turn is only possible when σ^2 is known a priori. When σ^2 is unknown, we need a formula in terms of its estimated value, $\hat{\sigma}^2$, leading to

$$\frac{c^T \hat{\beta} - c^T \beta}{\|c\|_{\hat{\mathcal{I}}}} \sim t_{n-p},$$

where $\hat{\mathcal{I}} = V^T V / \hat{\sigma}^2$ is the estimated Fisher information matrix. Thus, when σ^2 is unknown, a $100(1 - \alpha)\%$ confidence interval on $c^T \beta$ has endpoints

$$c^T \hat{\beta} \pm t_{n-p}^{1-\alpha/2} \cdot \|c\|_{\hat{\mathcal{I}}}, \quad (5)$$

where $t_{n-p}^{1-\alpha/2}$ is the upper $\alpha/2$ quantile of a Student's t distribution with $n - p$ degrees of freedom.

These facts can be used for testing $H_0 : c^T \beta = d$ vs. the alternative, $H_1 : c^T \beta \neq d$ for particular values of $c \in \mathbf{R}^p$ and d . Under the null hypothesis,

$$T_c := \frac{c^T \hat{\beta} - d}{\|c\|_{\hat{\mathcal{I}}}} \sim t_{n-p}, \quad (6)$$

so the p-value associated with the test is simply $\Phi_{t_{n-p}}(-|T_c|) + (1 - \Phi_{t_{n-p}}(|T_c|))$, where $\Phi_{t_{n-p}}$ is the cumulative distribution function for a Student's t distribution with $n - p$ degrees of freedom. Under a particular alternative hypothesis, say $c^T \beta = d'$, T_c is distributed as a noncentral Student's t with $n - p$ degrees of freedom and noncentrality parameter $\lambda = (d' - d)/\|c\|_{\mathcal{I}}$. I could not find a derivation of this, so here is one:

$$T_c = \frac{\frac{c^T \hat{\beta} - d'}{\|c\|_{\mathcal{I}}} + \frac{d' - d}{\|c\|_{\mathcal{I}}}}{\|c\|_{\hat{\mathcal{I}}} / \|c\|_{\mathcal{I}}} = \frac{\frac{c^T \hat{\beta} - d'}{\|c\|_{\mathcal{I}}} + \frac{d' - d}{\|c\|_{\mathcal{I}}}}{\sqrt{\hat{\sigma}^2 / \sigma^2}} = \frac{\frac{c^T \hat{\beta} - d'}{\|c\|_{\mathcal{I}}} + \frac{d' - d}{\|c\|_{\mathcal{I}}}}{\sqrt{\frac{(n-p) \cdot \hat{\sigma}^2 / \sigma^2}{n-p}}}.$$

From Equation (3), the first term in the numerator has a standard Gaussian distribution. The second term in the numerator is the noncentrality parameter. The denominator is the square root of a χ_{n-p}^2 random variable divided by its degrees of freedom. The numerator is a function of $\hat{\beta}$ while the denominator is a function of $\hat{\sigma}^2$, so the numerator and denominator are statistically independent. This is precisely the characterization of a noncentral Student's t distribution. For a test of size α , we would reject the null if $|T_c| > t_{n-p}^{1-\alpha/2}$. The probability of doing so under a particular alternative hypothesis is the power of the test and is given by:

$$1 - \Phi_{t_{n-p}; \lambda}(t_{n-p}^{1-\alpha/2}) + \Phi_{t_{n-p}; \lambda}(-t_{n-p}^{1-\alpha/2}),$$

where $\Phi_{t_{n-p}; \lambda}$ is the cumulative distribution function for a noncentral Student's t distribution with $n - p$ degrees of freedom and noncentrality parameter $\lambda = (d' - d)/\|c\|_{\mathcal{I}}$. Note we need to assume a value of σ^2 associated with the alternative hypothesis to compute $\|c\|_{\mathcal{I}}$ for the noncentrality parameter.

As special cases of the above discussion, when $c = \hat{e}_i$ (that is, a vector with a 1 in the i th entry, and zeros elsewhere), we get confidence intervals for, and hypothesis tests regarding, β_i . When $c = \nu(x_{\text{new}})$, we get confidence intervals for $\mu(x_{\text{new}})$; that is, the mean response of the model applied to a new data point. Confidence intervals on $Y \mid X = x_{\text{new}}$ involve an extra component of uncertainty due to the variance of the conditional distribution: even if we knew $\mu(x_{\text{new}})$ perfectly, the conditional distribution still has variance σ^2 . This leads to an extra term of $\hat{\sigma}^2$ as compared to Equation (5):

$$\nu(x_{\text{new}})^T \hat{\beta} \pm t_{n-p}^{1-\alpha/2} \cdot \sqrt{\hat{\sigma}^2 + \|\nu(x_{\text{new}})\|_{\hat{\mathcal{I}}}^2},$$

which are the endpoints of a confidence interval on Y [Wei05, §3.6].

When we want simultaneous confidence intervals on multiple linear combinations of β , we must adjust for the multiple comparisons. There is more than one approach to doing so. Suppose we are interested in $C\beta$, where $C \in \mathbf{R}^{q \times p}$. Then

$$C\hat{\beta} \sim \mathcal{N}(C\beta, C\hat{\mathcal{I}}^{-1}C^T),$$

which shows that $C\hat{\beta}$ is an unbiased estimator for $C\beta$, and that $C\hat{\beta}$ is normally distributed. Note that when $C = V$, we get the distribution of the fitted means, $\hat{\mu}$, since $\hat{\mu} := V\hat{\beta}$. Suppose we are interested in simultaneous confidence intervals on $C\beta$ at level α . The Bonferroni correction defines $\alpha' = \alpha/q$ and then simply uses the endpoints in Equation (5) for each individual component of $C\beta$, substituting α' for α . For example, suppose we wanted a simultaneous 95% confidence interval on $C\beta$, where C has $q = 5$ rows. Then $\alpha = 0.05$ and $\alpha' = 0.01$. So we would compute 99% confidence intervals on each component $c^{(i)}\beta$, where $c^{(i)}$ is the i th row of C . This approach is simple but overly conservative in many cases [Wei05, § 9.1.3].

Another method, due to Scheffé, defines simultaneous confidence intervals for *any* linear function of $C\beta$ [Sch59]. This is especially helpful when q is very large (in that case, the Bonferroni correction renders the confidence intervals too wide to be practically useful). Suppose C has rank r , and that $A \in \mathbf{R}^{r \times p}$ is any collection of r linearly independent rows of C . We wish to estimate simultaneous confidence intervals on quantities of the form $h^T A\beta$.² For example, when $h = \hat{e}_i$, this is simply the i th component of $A\beta$, but h can be anything. Then

$$h^T A\hat{\beta} \pm (r \cdot F_{r, n-p}^{1-\alpha/2})^{1/2} \cdot \|A^T h\|_{\hat{\mathcal{I}}}$$

is a $100(1 - \alpha)\%$ confidence interval on $h^T A\beta$, where $F_{r, n-p}^{1-\alpha/2}$ is the upper $\alpha/2$ quantile of Snedecor's F distribution having r and $n - p$ degrees of freedom in the numerator and denominator, respectively [SL03, § 5.1.1].

Now suppose C is full rank, and that $q < p$. We would like to test the hypothesis $C\beta = d$. As derived in [Woo17, § 1.3.4],

$$T_C := \frac{1}{q}(C\hat{\beta} - d)^T(C\hat{\mathcal{I}}^{-1}C^T)^{-1}(C\hat{\beta} - d) = \frac{1}{q}\|C\hat{\beta} - d\|_{C\hat{\mathcal{I}}^{-1}C^T}^2 \sim F_{q, n-p},$$

²Since the range of A^T is equal to the range of C^T , for any vector $h' \in \mathbf{R}^q$, there exists $h \in \mathbf{R}^r$ such that $C^T h' = A^T h$.

where $F_{q,n-p}$ is Snedecor's F distribution. This relationship generalizes (6) since an F distribution with one degree of freedom in the numerator is equivalent to a t^2 distribution [Wei05, § 3.5.3]. The p-value for this test would be $1 - \Phi_{F_{q,n-p}}(T_C)$, where $\Phi_{F_{q,n-p}}$ is the cumulative distribution function for an F distribution with the specified degrees of freedom.

Under a particular alternative hypothesis, say $C\hat{\beta} = d'$, T_C has a noncentral F distribution with noncentrality parameter $\lambda = \|d' - d\|_{C\mathcal{I}^{-1}C^T}^2$. As above, I cannot find the derivation of this anywhere, so I'll provide it here. Let $L^T L = (C\mathcal{I}^{-1}C^T)^{-1}$ (for example, L is the Cholesky decomposition). Then $L(C\hat{\beta} - d') \sim \mathcal{N}(0, I)$ and $L(C\hat{\beta} - d) \sim \mathcal{N}(L(d' - d), I)$. Under the alternative hypothesis,

$$\begin{aligned} T_C &= \frac{1}{q} (C\hat{\beta} - d)^T (C\hat{\mathcal{I}}^{-1}C^T)^{-1} (C\hat{\beta} - d) \\ &= \frac{(C\hat{\beta} - d)^T (C\mathcal{I}^{-1}C^T)^{-1} (C\hat{\beta} - d)}{\frac{q}{(n-p)\hat{\sigma}^2/\sigma^2}} \\ &= \frac{\frac{\|L(C\hat{\beta} - d)\|_2^2}{q}}{\frac{(n-p)\hat{\sigma}^2/\sigma^2}{n-p}}. \end{aligned}$$

The denominator is a χ_{n-p}^2 random variable divided by its degrees of freedom. The numerator is the sum of squares of independent unit variance Gaussian random variables with mean vector $\mu = L(d' - d)$, so letting

$$\begin{aligned} \lambda &= \mu^T \mu \\ &= (d' - d)^T L^T L (d' - d) \\ &= (d' - d)^T (C\mathcal{I}^{-1}C^T)^{-1} (d' - d) \\ &= \|d' - d\|_{C\mathcal{I}^{-1}C^T}^2, \end{aligned}$$

we see that the numerator is a noncentral χ_q^2 random variable with noncentrality parameter λ , divided by its degrees of freedom. Since the numerator is a function of $\hat{\beta}$, and the denominator is a function of $\hat{\sigma}^2$, the numerator and denominator are statistically independent, which demonstrates the test statistic is F -distributed. For a test of size α , we would reject the null if $T_C > F_{q,n-p}^{1-\alpha}$. The probability of doing so under a particular hypothesis is the power of the test and is given by:

$$1 - \Phi_{F_{q,n-p;\lambda}}(F_{q,n-p}^{1-\alpha}),$$

where $\Phi_{F_{q,n-p;\lambda}}$ is the cumulative distribution function of a noncentral F distribution with the stated degrees of freedom and noncentrality parameter.

A.3 Model Selection

The F -test is useful in several situations. Suppose we are considering a sequence of nested models: a model consisting only of a grand mean (in which μ does not depend on the features

at all), a model consisting only of main effects, a model with first-order interactions, and a model with higher-order interactions. We may wish to check H_0 : the model consists only of main effects vs. H_1 : interactions are present. This amounts to checking whether $C\beta = 0$, where $C\beta$ are the components of β corresponding to the interaction terms. Or we may want to check $H_0 : \mu(x) = \mu$ vs. $H_1 : \mu(x) \neq \mu$, where μ is the grand mean. In this case, we are checking whether there is any evidence of the mean response depending on the features. Of course, checking multiple hypotheses is subject to the multiple comparisons problem discussed above.

If a particular feature is categorical with at least three levels, it will consist of at least two parameters. We would typically want to test whether all associated parameters are non-zero, not just one. Or if $\nu(x)$ consists of multiple transformations of a particular feature, like $\nu(x) = [\cdots \ x_1 \ x_1^2 \ \log(x_1) \ \cdots]^T$, we might want to simultaneously test all the corresponding components of β for being non-zero. The F -test described here is applicable to these scenarios.

The F -test is not the only approach to model selection, however; other methods include those based on information criteria like Akaike's Information Criterion (AIC) [SL03, § 12.3.3], the Bayesian Information Criterion (BIC), and Mallows's C_p statistic [Wei05, § 10.2.1]. These are given, respectively, by:

$$\begin{aligned} \text{AIC} &= n \log(2\pi\sigma^2) + \frac{1}{\sigma^2} \|y - V\hat{\beta}\|_2^2 + 2 \cdot \text{dof} \quad (\sigma^2 \text{ known}) \\ \text{AIC} &= n \log(2\pi \|y - V\hat{\beta}\|_2^2 / n) + n + 2 \cdot \text{dof} \quad (\sigma^2 \text{ unknown}) \\ \text{AICc} &= \text{AIC} + \frac{2\text{dof}^2 + 2\text{dof}}{n - \text{dof} - 1} \quad (\text{linear model correct}) \\ &= n \log(2\pi \|y - V\hat{\beta}\|_2^2 / n) + \frac{n(n + \text{dof} - 1)}{n - \text{dof} - 1} \quad (\sigma^2 \text{ unknown}) \\ \text{BIC} &= n \log(\|y - V\hat{\beta}\|_2^2 / n) + \text{dof} \cdot \log(n) \\ C_p &= \frac{\|y - V\hat{\beta}\|_2^2}{\sigma^2} + 2 \cdot \text{dof} - n, \end{aligned}$$

where dof is equal to the number of parameters estimated in the model. If σ^2 is known a priori, this is simply p ; otherwise, it is $p + 1$. Note that AIC comes in two forms depending on whether σ^2 is known a priori. A modification of AIC is often desirable for small sample sizes; this is known as the corrected AIC, or AICc. The correction term depends on whether we believe the model truly is normally distributed with a mean depending linearly on the parameters; this is the formula shown [BA02, § 7.7.6].

These formulae clearly illustrate the tradeoff between a better fitting model and a model having more parameters. Notably, the BIC formula penalizes degrees of freedom much more strongly than does the AIC, and thus will lead to simpler models. Caution is advised when using C_p with unknown σ^2 [Woo17, § 1.8.6]. If we must, it is best to estimate σ^2 using the most flexible model under consideration (that is, the model with all parameters included), and using the same value for all models being compared.

Cross validation is another, more computationally intensive approach to model selection. By dividing the data set into a training and test set, we fit the model to the training set and use the result to predict the response for the data in the test set, using the actual response to compute the prediction error. The model giving the best prediction error is the one we select. Often we will then refit the model on the entire data set [HTF01, § 7.10].

A.4 Checking Model Assumptions

Finally, we discuss the model residuals, $\hat{\epsilon}^{(i)} = y^{(i)} - \hat{\mu}(x^{(i)})$. We have already been using the residuals to estimate the variance, σ^2 , but examining the residuals is also useful for investigating departures from the assumptions of the linear model: that the response is normally distributed, that the mean of this distribution depends linearly on the model parameters, that the variance is constant (or is of the form $\sigma^2/w^{(i)}$ with known $w^{(i)}$), and that the observations are statistically independent. These assumptions may be checked by graphing the residuals.

Let $H = V(V^T V)^{-1} V^T$ be the so-called *hat matrix*. Then $\hat{\epsilon} \sim \mathcal{N}(0, \sigma^2(I - H))$. Notably, the residuals are correlated, and have different variances (however, the residuals are statistically independent of $\hat{\mu}$). Because of this, it is typical to standardize the residuals so that they have equal variance. The internally and externally Studentized residuals are defined as

$$r^{(i)} = \frac{\hat{\epsilon}^{(i)}}{\sqrt{\hat{\sigma}^2 \cdot (1 - h_i)}}$$

$$t^{(i)} = \frac{\hat{\epsilon}^{(i)}}{\sqrt{\hat{\sigma}_{(i)}^2 \cdot (1 - h_i)}},$$

respectively, where h_i is the i th diagonal element of H and $\hat{\sigma}_{(i)}^2 = \frac{1}{n-p-1} \sum_{j \neq i} \hat{\epsilon}^{(j)2}$. As [SL03, § 10.2] states, $(r^{(i)})^2/(n-p)$ has a $\text{beta}[\frac{1}{2}, \frac{1}{2}(n-p-1)]$ distribution which means they are identically distributed (but not independent). The externally Studentized residuals, $t^{(i)}$, have a t_{n-p-1} distribution. The externally Studentized residuals are less prone to outliers than are the internally Studentized residuals.

Consider a graph of $t^{(i)}$ (or $r^{(i)}$) against $\hat{\mu}(x^{(i)})$. If the model is correct, we expect to see a scattering of points with no discernible pattern, since the Studentized residuals would be identically distributed and independent of the mean response. If there is an apparent trend in the residuals, that may indicate a nonlinearity in the model.

To assess a potential dependence between the mean response and the variance, [SL03, § 10.4.2] recommends plotting the squared residuals, $\epsilon^{(i)2}$, against the fitted means, $\hat{\mu}(x^{(i)})$. If the variance increases with the mean response, this plot will exhibit a wedge shape. We can apply a smoother such as lowess to estimate the relationship between the mean response and the variance. This gives an estimate of the variance associated with each observation, which can then be used to determine weights for the observations. Since the variance of the i th observation is assumed to be $\sigma^2/w^{(i)}$, and the estimated variance of the i th observation is $(\epsilon^{(i)})^2$, we have $w^{(i)} = (\epsilon^{(i)})^{-2}$, where we are setting $\sigma^2 = 1$ since we are

directly estimating the variance of each individual observation. Iterating on this procedure (estimating the model using weighted least squares, plotting the squared residuals against the mean response, smoothing this plot to estimate the variance of each observation) gives an estimate of β that is asymptotically as efficient as knowing the weights a priori.

We can test the normality assumption using a Q-Q plot, which graphs the observed quantiles of the raw residuals against the quantiles of a standard Gaussian distribution. Alternatively we could graph the quantiles of the Studentized residuals against the quantiles of their theoretical distributions [SL03, § 10.5.1].

One of the assumptions of the linear model is that the observations are independent. If the observations have a known correlation structure, various approaches exist for fitting models. If we believe the observations are independent, we can check for a specific deviation from this assumption called *serial correlation*. That is, we can check for correlations between sequential pairs of observations. This is especially relevant when the order of observations is physically meaningful. In the absence of correlation, a residual with positive sign is equally likely to be followed by a residual with positive or negative sign, which can easily be examined graphically. A significance test-based procedure was discussed in a series of papers by Durbin and Watson. This test checks for a first-order autoregressive model for the residuals: $\hat{\epsilon}^{(i)} = \rho\hat{\epsilon}^{(i-1)} + \delta^{(i)}$, where $\delta^{(i)}$ are independent normal variables. Let

$$D = \frac{\sum_{i=2}^n (\hat{\epsilon}^{(i)} - \hat{\epsilon}^{(i-1)})^2}{\sum_{i=1}^n (\hat{\epsilon}^{(i)})^2} = \frac{\hat{\epsilon}^T A \hat{\epsilon}}{\hat{\epsilon}^T \hat{\epsilon}}, \text{ where}$$

$$A = \begin{bmatrix} 1 & -1 & 0 & 0 & \cdots & \cdots & 0 \\ -1 & 2 & -1 & 0 & \cdots & \cdots & \cdots \\ 0 & -1 & 2 & -1 & \cdots & \cdots & \cdots \\ 0 & 0 & -1 & 2 & \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots & \cdots & \cdots & 2 & -1 \\ \cdots & \cdots & \cdots & \cdots & \cdots & -1 & 1 \end{bmatrix}$$

Under the null hypothesis of independent observations, D has the same distribution as

$$r = \frac{\sum_{i=1}^{n-p} \zeta_i \zeta_i^2}{\sum_{i=1}^{n-p} \zeta_i^2}, \quad (7)$$

where ζ_i are IID standard Gaussian variables and ξ_i are the nonzero eigenvalues of $(I - H)A$ —assuming V is full rank, there will be exactly $n - p$ of these [DW50, pg. 416]. Two issues present themselves, one computational, the other theoretical. Computing the eigenvalues of $(I - H)A$ may be computationally intensive if $n - p$ is gigantic or if you happen to be living in 1950. More problematically, exact tail probabilities for distributions of the form (7) are not available.

Since H depends explicitly on the features, it will be different for each regression analysis; however, it can be shown that the eigenvalues ξ_i of $(I - H)A$ are bounded by pairs of the eigenvalues of A , λ_i [DW50]. Specifically, if we sort the eigenvalues so that $\xi_1 \leq \xi_2 \leq \cdots \xi_{n-p}$

and $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$, then $\lambda_i \leq \xi_i \leq \lambda_{i+p}$, $i = 1, 2, \dots, (n-p)$. The eigenvalues of A have a simple form: $\lambda_j = 2 \cdot (1 - \cos(\pi(j-1)/n))$, $j = 1, 2, \dots, n$ [DW50, pg. 426].

Tighter bounds hold whenever s of the eigenvectors of $(I - H)A$ are linear combinations of s of the eigenvectors of A . That doesn't seem like it would happen very often, but since $\mathbf{1}$ is an eigenvector of A (corresponding to an eigenvalue of zero), when the model includes a constant affine term, then one of the columns of V is $\mathbf{1}$, and $s \geq 1$. When that happens, we may discard the corresponding s eigenvalues of A , leaving $n - s$ eigenvalues λ_i , and the bounds become $\lambda_i \leq \xi_i \leq \lambda_{i+p-s}$.

These bounds on the eigenvalues of $(I - H)A$ become bounds on the distribution of the Durbin-Watson statistic: $r_L \leq r \leq r_U$, where

$$r_L = \frac{\sum_{i=1}^{n-p} \lambda_i \zeta_i^2}{\sum_{i=1}^{n-p} \zeta_i^2},$$

$$r_U = \frac{\sum_{i=1}^{n-p} \lambda_{i+p-s} \zeta_i^2}{\sum_{i=1}^{n-p} \zeta_i^2}.$$

Note that the distributions of r_L and r_U depend on n , p , and s , but not on the features.

It is straightforward to show that $r_L \geq \lambda_1 = 0$ and $r_U \leq \lambda_n < 4$, which shows that the Durbin-Watson statistic satisfies $0 \leq D < 4$. In the presence of positive serial correlation, $\rho > 0$, D will tend to be closer to 0. When $\rho < 0$, D will tend to be closer to 4. In the absence of serial correlation, D will typically be close to 2. Let Φ_L , Φ_U , and Φ_{DW} be the cumulative distribution functions of r_L , r_U , and r , respectively, so that, for example, $\Phi_L^{-1}(d) = \mathbf{Prob}\{r_L \leq d\}$. In light of the above discussion, $\Phi_L^{-1}(d) \geq \Phi_{DW}^{-1}(d) \geq \Phi_U^{-1}(d)$ [DW50, pg. 418]. Table 1 shows how these functions provide bounds on the p-values for various tests related to serial correlation. Evaluating the exact p-values require the eigenvalues ξ_i , which depend on the features. Evaluating the bounds only requires tail-probabilities for Φ_L and Φ_U , which do not depend on the features (but do depend on n , p , and s). These have been tabulated for various values of n , p , and s , for example in [DW51].

When performing an analysis with numbers of observations and parameters not represented in an available table, we are still left with the challenge of computing the tail probabilities of r_L and r_U . We may proceed by approximating $r_L/4$ and $r_U/4$ as beta-distributed. Tail probabilities for the beta distribution may then be mapped to p-values for the Durbin-Watson statistic. The beta distributions are chosen to have the same means and variances as $r_L/4$ and $r_U/4$, respectively. These may be expressed in terms of the eigenvalues of A , λ_i .

Test	p-value	Lower bound	Upper bound
$\rho = 0$ vs. $\rho > 0$	$\Phi_{DW}^{-1}(d)$	$\Phi_U^{-1}(d)$	$\Phi_L^{-1}(d)$
$\rho = 0$ vs. $\rho < 0$	$1 - \Phi_{DW}^{-1}(d)$	$1 - \Phi_L^{-1}(d)$	$1 - \Phi_U^{-1}(d)$
$\rho = 0$ vs. $\rho \neq 0$	$\Phi_{DW}^{-1}(d') + 1 - \Phi_{DW}^{-1}(d')$	$\Phi_U^{-1}(d') + 1 - \Phi_L^{-1}(d')$	$\Phi_L^{-1}(d') + 1 - \Phi_U^{-1}(d')$

Table 1: Bounds on p-values for one and two-sided tests regarding the correlation parameter, ρ . In all cases, d is the observed value of the Durbin-Watson statistic. In the last row, $d' = 2 - |2 - d|$.

Statistic	Mean	Variance
r	$\mu = \frac{1}{n-p} \sum_{i=1}^{n-p} \xi_i$	$\sigma^2 = \frac{2 \sum_{i=1}^{n-p} (\xi_i - \mu)^2}{(n-p)(n-p+2)}$
r_L	$\mu_L = \frac{1}{n-p} \sum_{i=1}^{n-p} \lambda_i$	$\sigma_L^2 = \frac{2 \sum_{i=1}^{n-p} (\lambda_i - \mu_L)^2}{(n-p)(n-p+2)}$
r_U	$\mu_U = \frac{1}{n-p} \sum_{i=1}^{n-p} \lambda_{i+p-s}$	$\sigma_U^2 = \frac{2 \sum_{i=1}^{n-p} (\lambda_{i+p-s} - \mu_U)^2}{(n-p)(n-p+2)}$

Table 2: Means and variances of the Durbin-Watson and related statistics

If the eigenvalues, ξ_i , are in fact known, we can dispense with the bounds entirely. Means and variances for all three statistics are reported in Table 2.

A beta distribution is typically characterized by parameters α and β . The mean of a beta distribution is $\frac{\alpha}{\alpha+\beta}$ and the variance is $\frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$, so the beta distribution used to approximate r , for example, has parameters

$$\alpha = \frac{\mu^2(4-\mu)}{4\sigma^2} - \frac{\mu}{4},$$

$$\beta = \frac{\mu(4-\mu)^2}{4\sigma^2} - \frac{4-\mu}{4}.$$

For the purposes of calculating p-values, $\Phi_{\alpha,\beta}^{-1}(d/4)$ may be substituted for $\Phi_{\text{DW}}^{-1}(d)$ in Table 1, where $\Phi_{\alpha,\beta}$ is the cumulative distribution function of a beta random variable with parameters α and β . For more details and discussion of alternative approaches, see [DW71].

Finally, we want to examine how any potential outliers affect the fitted model. Small changes to the response for observations on the outskirts of the feature space can have a big effect on the model; such points are said to have high *leverage*. Leverage may be investigated by examining the diagonal terms of the hat matrix, h_i . The higher a particular h_i , the larger the influence of the corresponding observation on the fitted model. We might consider any observation having $h_i > 2p/n$ to have high leverage [SL03, 10.6.1]. An observation that does not have high leverage, but deviates wildly from the mean response can also have undue influence on the model. Since the externally Studentized residuals have a t_{n-p-1} distribution, any residual with $|t^{(i)}| > 2$ should be examined (corresponding approximately to the upper and lower 2.5% quantiles).

If a particular observation is both an outlier and has high leverage, we can try omitting the observation, or reducing its weight, and refitting. Large changes to the fitted model indicate the point has high influence. Deciding whether or not to remove the observation depends on the goals of the analysis, how the data were collected, and so forth. A variety of statistics are available for quantifying the impact of leaving out a single observation without actually having to refit the model[SL03, § 10.6.3]. For example, the impact of leaving out the i th observation to the i th fitted value is $h_i \epsilon^{(i)} / (1 - h_i)$. This can be standardized giving $t^{(i)} \sqrt{h_i / (1 - h_i)}$. A cutoff of $2\sqrt{p/(n-p)}$ can be used for identifying high influence points. Another statistic, called Cook's D , may be written:

$$D^{(i)} = \left(r^{(i)}\right)^2 \frac{h_i}{p(1 - h_i)}.$$

Cook recommended using $F_{p,n-p}^{0.10}$ as the cutoff for identifying high influence points [Coo77].

When one or more observations have been identified as possible outliers, a formal test may be applied [SL03, § 10.6.4]. If we are testing a set of k observations, we augment ν with k extra entries. The j th of these entries is 1 for the j th potential outlier, and zero otherwise. We fit the expanded model. Let $\hat{\gamma}$ be the subset of entries of $\hat{\beta}$ corresponding to these extra entries of ν . If the unaugmented model is correct, and the observations under consideration are *not* outliers, then $\gamma = 0$. The F -test outlined in § A.2 may be applied to test this hypothesis, giving a p-value for the collection of potential outliers. In this case, the test statistic is:

$$T_C = \frac{1}{k} \|\hat{\gamma}\|_{\hat{\mathcal{I}}_{\text{augmented}}^{-1}}^2 \sim F_{k,n-p-k},$$

where $\hat{\mathcal{I}}_{\text{augmented}}^{-1}$ is the $k \times k$ submatrix of $\hat{\mathcal{I}}^{-1}$ corresponding to the augmented entries of ν . The p-value is $1 - \Phi_{F_{k,n-p-k}}(T_C)$. If this p-value is small, that constitutes evidence that at least one of the observations under consideration is an outlier.

References

- [BA02] Kenneth P. Burnham and David Raymond Anderson. *Model Selection and Multimodel Inference, A Practical Information-Theoretic Approach*. Springer, 2nd edition, 2002.
- [BD87] George Edward Pelham Box and Norman Richard Draper. *Empirical Model-Building and Response Surfaces*. John Wiley & Sons, Inc., 1987.
- [BD93] David Spencer Birkes and Yadolah Dodge. *Alternative Methods of Regression*. John Wiley & Sons, Inc., 1993.
- [CB01] George C. Casella and Roger L. Berger. *Statistical Inference*. Cengage Learning, 2nd edition, 2001.
- [Coo77] Ralph Dennis Cook. Detection of influential observation in linear regression. *Technometrics*, 19(1):15–18, 1977.
- [DW50] James Durbin and Geoffrey Stuart Watson. Testing for serial correlation in least squares regression. I. *Biometrika*, 37(3-4):409–428, 1950.
- [DW51] James Durbin and Geoffrey Stuart Watson. Testing for serial correlation in least squares regression. II. *Biometrika*, 38(1-2):159–177, 1951.
- [DW71] James Durbin and Geoffrey Stuart Watson. Testing for serial correlation in least squares regression. III. *Biometrika*, 58(1):1–19, 1971.
- [ET93] Bradley Efron and Robert John Tibshirani. *An Introduction to the Bootstrap*. Chapman & Hall, 1993.
- [HT86] Trevor John Hastie and Robert John Tibshirani. Generalized additive models. *Statist. Sci.*, 1(3):297–310, 1986.
- [HTF01] Trevor John Hastie, Robert John Tibshirani, and Jerome Harold Friedman. *The Elements of Statistical Learning*. Springer, 2001.
- [MN89] Peter McCullagh and John Ashworth Nelder. *Generalized Linear Models*. Chapman & Hall, 2nd edition, 1989.
- [NW72] John Ashworth Nelder and Robert William MacLagan Wedderburn. Generalized linear models. *Journal of the Royal Statistical Society. Series A*, 135(3):370–384, 1972.
- [Sch59] Henry Scheffé. *The Analysis of Variance*. John Wiley & Sons, Inc., 1959.
- [SL03] George Arthur Frederick Seber and Alan James Lee. *Linear Regression Analysis*. John Wiley & Sons, Inc., 2nd edition, 2003.

- [Str12] Walter W. Stroup. *Generalized Linear Mixed Models*. Chapman & Hall, 2012.
- [Tib96] Robert John Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B*, 58(1):267–288, 1996.
- [Wei05] Sanford Weisberg. *Applied Linear Regression*. John Wiley & Sons, Inc., 3rd edition, 2005.
- [Woo17] Simon N. Wood. *Generalized Additive Models, An Introduction with R*. Chapman & Hall, 2nd edition, 2017.