

Lecture 2: Logistic Regression

STATS305B: Applied Statistics II

Scott Linderman

January 8, 2025

Recap

Last time...

- ▶ Basic distributions (Bernoulli, binomial, Poisson, categorical, multinomial)
- ▶ Basic inference (MLE, hypothesis tests, confidence intervals)
- ▶ Contingency tables (likelihood ratio test of independence)

Outline

Today, from contingency tables to logistic regression

- ▶ Setting up the model
- ▶ Parameter estimation
 - ▶ The hacky way: OLS
 - ▶ Maximum likelihood estimation
 - ▶ Gradient descent and Newton's method
 - ▶ Iteratively reweighted least squares (IRLS)
- ▶ Regularization
- ▶ Convergence rates

Contingency Tables with Binary Responses

A special case of contingency table analyses is when $X \in \{1, \dots, l\}$ corresponds to a categorical feature (e.g., to which of l groups a data point belongs) and $Y \in \{0, 1\}$ is a binary response.

The corresponding tables are $l \times 2$. Normalizing each row yields a Bernoulli conditional distribution,

$$Y \mid X = i \sim \text{Bern}(\pi_{1|i}).$$

where, $\pi_{1|i} = \Pr(Y = 1 \mid X = i)$ or, equivalently, $\pi_{1|i} = \mathbb{E}[Y \mid X = i]$.

Conditional distributions are often the primary objects of interest.

One, Two, Many...

What if we have more than one feature, X_1, \dots, X_p ?

What if the features take on continuous values?

Idea: Model the conditional distribution directly. Let

- ▶ $Y \in \{0, 1\}$ denote a binary response
- ▶ $\mathbf{X} \in \mathbb{R}^p$ denote associated covariates.

E.g., Y could denote whether or not your favorite football team wins their match, and X could represent features of the match like whether its a home or away game, who their opponent is, etc.

One, Two, Many...

Model the conditional distribution as,

$$Y \mid \mathbf{X} = \mathbf{x} \sim \text{Bern}(\pi(\mathbf{x}))$$

where

$$\pi(\mathbf{x}) = \Pr(Y = 1 \mid \mathbf{X} = \mathbf{x}) = \mathbb{E}[Y \mid \mathbf{X} = \mathbf{x}].$$

Standard regression setup.

Modeling choice: what functional form of $\pi(\mathbf{x})$?

Linear Regression

If you took STATS 305A, you know pretty much everything there is to know about linear regression with continuous response variables, $Y \in \mathbb{R}$. Why don't we just apply that same model to binary responses? Specifically, let,

$$\pi(\mathbf{x}) = \boldsymbol{\beta}^\top \mathbf{x} = \sum_{j=1}^p \beta_j x_j.$$

Problem: linear model produces probabilities or expectations $\pi(\mathbf{x})$ outside $[0, 1]$

Logistic Regression

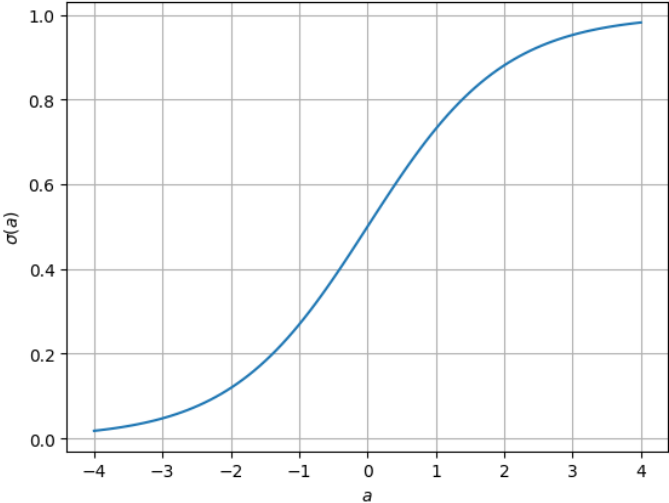
The idea is simple: keep the linear part of linear regression, but apply a **mean (aka inverse link) function**, $f : \mathbb{R} \mapsto [0, 1]$, to ensure $\pi(\mathbf{x})$ returns valid probabilities,

$$\pi(\mathbf{x}) = f(\boldsymbol{\beta}^\top \mathbf{x}).$$

There are infinitely many squashing nonlinearities that we could choose for f , but a particularly attractive choice is the **logistic (aka sigmoid) function**,

$$f(a) = \frac{e^a}{1 + e^a} = \frac{1}{1 + e^{-a}} \triangleq \sigma(a).$$

Logistic Function



Nice Properties of the Logistic Function

- ▶ Monotonically increasing
- ▶ Invertible: $\beta^\top \mathbf{x}$ correspond to the **log odds** of the binary response since the inverse of the sigmoid function is the **logit function**,

$$\beta^\top \mathbf{x} = \sigma^{-1}(\pi(\mathbf{x})) = \log \frac{\pi(\mathbf{x})}{1 - \pi(\mathbf{x})}.$$

- ▶ Simplifies some calculations for MLE

Another common mean function is the Gaussian CDF, and we'll consider that in a later chapter.

Hacky Parameter Estimation with OLS

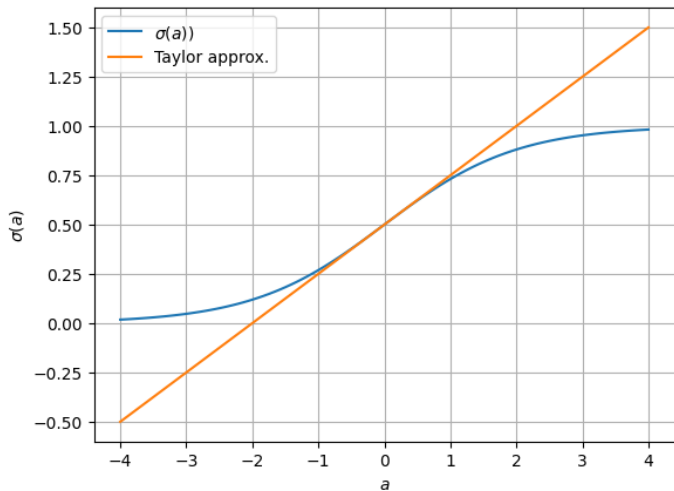
How do we estimate the parameters, $\hat{\beta}$?

If the model were linear, then our first inclination might be to use ordinary least squares (OLS). Of course, the sigmoidal function above renders the model nonlinear.

Idea: What if we just used a Taylor approximation around the origin,

$$\sigma(a) \approx \sigma(0) + \sigma'(0)a = \frac{1}{2} + \frac{a}{4}$$

Hacky Parameter Estimation with OLS



Hacky Parameter Estimation with OLS

Then under this model,

$$\mathbb{E}[Y | \mathbf{X} = \mathbf{x}] \approx \frac{1}{2} + \frac{\boldsymbol{\beta}^\top \mathbf{x}}{4}$$

or equivalently,

$$\mathbb{E}[Z | \mathbf{X} = \mathbf{x}] \approx \boldsymbol{\beta}^\top \mathbf{x}$$

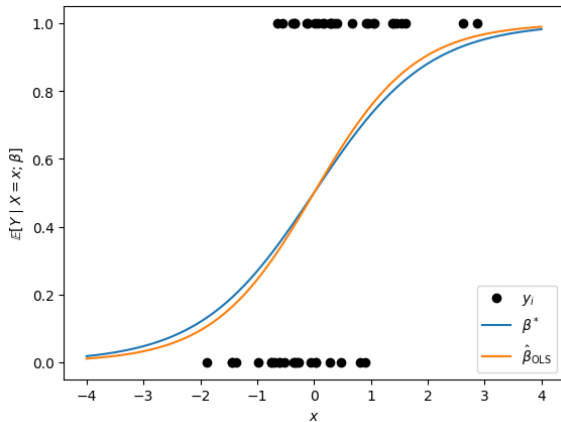
where $Z = 4(Y - \frac{1}{2}) \in \{-2, +2\}$ is an adjusted response variable.

Given a set of n observations of features and (adjusted) responses, $\{\mathbf{x}_i, z_i\}_{i=1}^n$, the OLS estimate is,

$$\hat{\boldsymbol{\beta}}_{\text{OLS}} = \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} \left(\sum_{i=1}^n \mathbf{x}_i z_i \right)$$

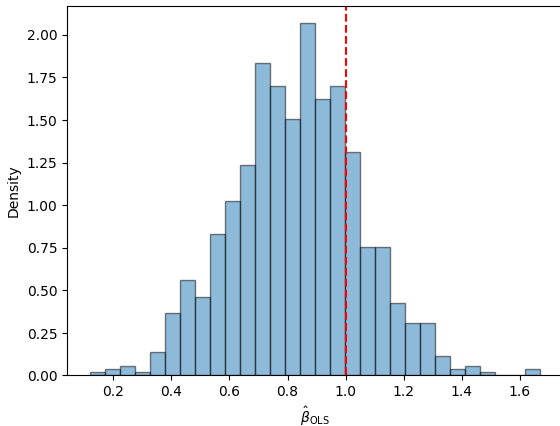
Demo

Simulated dataset with scalar covariates $X_i \sim N(0, 1)$ and responses $Y_i | X_i = x_i \sim \text{Bern}(\sigma(\beta^* x_i))$ for $\beta^* = 1$.



Demo

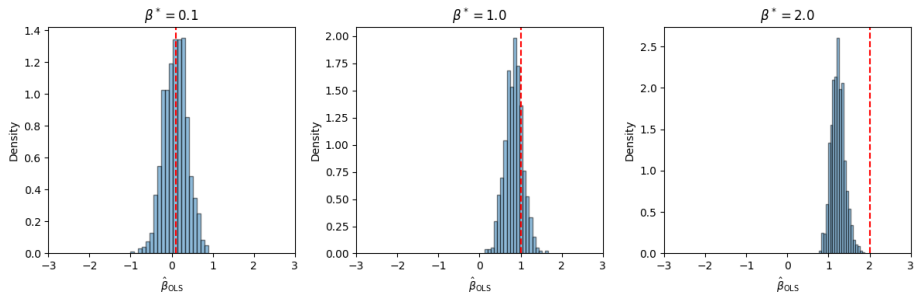
Not too shabby! Let's try it with a bunch of simulated datasets.



Demo

It looks a bit biased, but it's not terrible.

What if we try for other values of β^* ?



Question: Why do you think the OLS estimator becomes more biased as β^* grows?

Maximum Likelihood Estimation

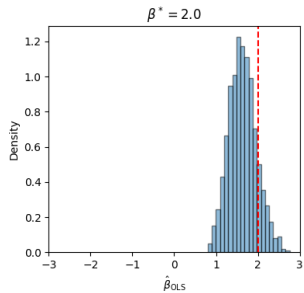
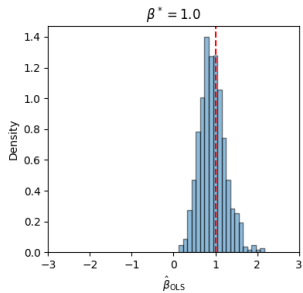
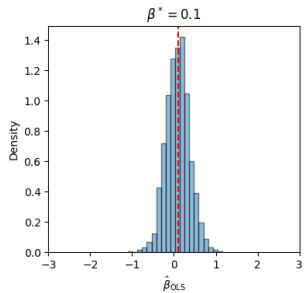
For standard linear regression with independent Gaussian responses and homoskedastic noise,

$$\hat{\beta}_{\text{MLE}} = \hat{\beta}_{\text{OLS}}.$$

Unfortunately, for logistic regression, there isn't a closed form for $\hat{\beta}_{\text{MLE}}$. However, we can use standard optimization techniques to compute it.

There are plenty of standard implementations of maximum likelihood estimation for logistic regression models. Let's see how scikit-learn's implementation fares on the simulated datasets above.

Demo



Deriving the MLE

Maximizing the likelihood is equivalent to minimizing the (average) negative log likelihood for a collection of covariates and responses, $\{\mathbf{x}_i, y_i\}_{i=1}^n$,

$$\begin{aligned}\mathcal{L}(\boldsymbol{\beta}) &= -\frac{1}{n} \sum_{i=1}^n \log \text{Bern}(y_i; \pi(\mathbf{x}_i)) \\ &= -\frac{1}{n} \sum_{i=1}^n y_i \log \pi(\mathbf{x}_i) + (1 - y_i) \log(1 - \pi(\mathbf{x}_i)) \\ &= -\frac{1}{n} \sum_{i=1}^n y_i \log \frac{\pi(\mathbf{x}_i)}{1 - \pi(\mathbf{x}_i)} + \log(1 - \pi(\mathbf{x}_i)) \\ &= -\frac{1}{n} \sum_{i=1}^n y_i \boldsymbol{\beta}^\top \mathbf{x}_i + \log(1 - \pi(\mathbf{x}_i)).\end{aligned}$$

Since $\boldsymbol{\beta}^\top \mathbf{x}_i$ is the log odds, as shown above.

Deriving the MLE

Plugging in the definition of $\pi(\mathbf{x}_i)$, the second term simplifies as well,

$$\begin{aligned}\mathcal{L}(\boldsymbol{\beta}) &= -\frac{1}{n} \sum_{i=1}^n \left[y_i \boldsymbol{\beta}^\top \mathbf{x}_i + \log \left(1 - \frac{e^{\boldsymbol{\beta}^\top \mathbf{x}_i}}{1 + e^{\boldsymbol{\beta}^\top \mathbf{x}_i}} \right) \right] \\ &= -\frac{1}{n} \sum_{i=1}^n \left[y_i \boldsymbol{\beta}^\top \mathbf{x}_i - \log \left(1 + e^{\boldsymbol{\beta}^\top \mathbf{x}_i} \right) \right].\end{aligned}$$

Computing the Gradient

To minimize the negative log likelihood, take the gradient,

$$\begin{aligned}\nabla \mathcal{L}(\boldsymbol{\beta}) &= -\frac{1}{n} \sum_{i=1}^n \left(y_i \mathbf{x}_i - \frac{e^{\boldsymbol{\beta}^\top \mathbf{x}_i}}{1 + e^{\boldsymbol{\beta}^\top \mathbf{x}_i}} \mathbf{x}_i \right) \\ &= -\frac{1}{n} \sum_{i=1}^n (y_i - \sigma(\boldsymbol{\beta}^\top \mathbf{x}_i)) \mathbf{x}_i.\end{aligned}$$

The gradient is a weighted sum of the covariates, and the weights are the residuals $y_i - \sigma(\mathbf{x}_i^\top \boldsymbol{\beta})$, i.e., the difference between the observed and expected response.

Intuition: move in the direction of covariates where the residual is positive (we are underestimating the mean), and opposite covariates where the residual is negative (where we are overestimating the mean).

Computing the Gradient

Algorithm:

$$\beta_{t+1} \leftarrow \beta_t - \alpha_t \nabla \mathcal{L}(\beta_t),$$

where $\alpha_t \in \mathbb{R}_+$ is the step-size at iteration i of the algorithm.

If the step sizes are chosen appropriately and the objective is well behaved, the algorithm converges to at least a local optimum of the log likelihood.

Convexity of the Log Likelihood

If the objective is convex, then all local optima are also global optima, and we can give stronger guarantees on gradient descent. To check the convexity of the log likelihood, we need to compute its Hessian,

$$\nabla^2 \mathcal{L}(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \sigma'(\boldsymbol{\beta}^\top \mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top$$

where $\sigma'(a)$ is the derivative of the logistic function. That is,

$$\begin{aligned} \sigma'(a) &= \frac{d}{da} \sigma(a) \\ &= \frac{e^a}{(1 + e^a)^2} \\ &= \sigma(a)(1 - \sigma(a)). \end{aligned}$$

Convexity of the Log Likelihood

Plugging this in,

$$\begin{aligned}\nabla^2 \mathcal{L}(\boldsymbol{\beta}) &= \frac{1}{n} \sum_{i=1}^n \sigma(\boldsymbol{\beta}^\top \mathbf{x}_i)(1 - \sigma(\boldsymbol{\beta}^\top \mathbf{x}_i)) \mathbf{x}_i \mathbf{x}_i^\top \\ &= \frac{1}{n} \sum_{i=1}^n w_i \mathbf{x}_i \mathbf{x}_i^\top,\end{aligned}$$

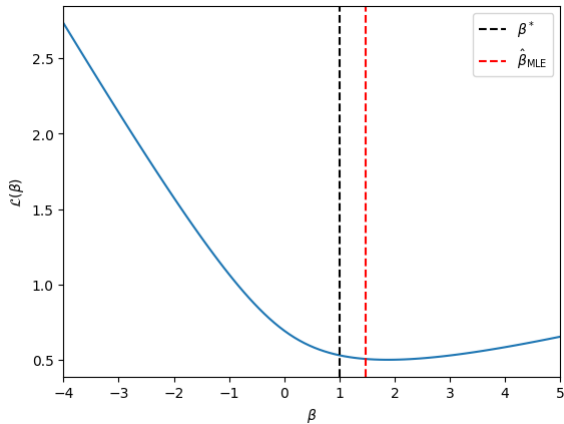
where the weights are, $w_i = \sigma(\boldsymbol{\beta}^\top \mathbf{x}_i)(1 - \sigma(\boldsymbol{\beta}^\top \mathbf{x}_i)) = \text{Var}[Y | \mathbf{X} = \mathbf{x}_i]$.

The Hessian is a weighted sum of outer products of covariates where the weights are equal to the conditional variance, which are non-negative.

Since $w_i > 0$, the Hessian is **positive semi-definite**, which implies that the negative log likelihood is convex.

Example

As an example, let's plot the negative log likelihood for a scalar covariate example, as above.



Pathologies in the Separable Regime

Suppose the two classes are *linearly separable*,

$$\exists \mathbf{u} \in \mathbb{S}_{p-1} \text{ s.t. } \begin{cases} \mathbf{u}^\top \mathbf{x}_i > 0 & \text{if } y_i = 1 \\ \mathbf{u}^\top \mathbf{x}_i < 0 & \text{if } y_i = 0 \end{cases}$$

Now let $\boldsymbol{\beta} = c\mathbf{u}$ for any $c \in \mathbb{R}_+$. We have

$$\lim_{c \rightarrow \infty} \sigma(c\mathbf{u}^\top \mathbf{x}_i) = y_i.$$

In this limit, the model is saturated: $\Pr(y_i | \mathbf{x}_i) = 1$ for all $i = 1, \dots, n$; the negative log likelihood goes to $\lim_{c \rightarrow \infty} \mathcal{L}(c\mathbf{u}) = 0$; and the MLE does not exist since $\boldsymbol{\beta}^*$ diverges.

If we run gradient descent in this setting, the magnitude of the estimate $\hat{\boldsymbol{\beta}}$ will grow without bound, while the objective converges to zero.

L_2 -Regularization

These pathologies can be averted with a little regularization,

$$\mathcal{L}(\boldsymbol{\beta}) = -\frac{1}{n} \sum_{i=1}^n \log \text{Bern}(y_i; \sigma(\boldsymbol{\beta}^\top \mathbf{x}_i)) + \frac{\gamma}{2} \|\boldsymbol{\beta}\|_2^2.$$

The regularizer penalizes larger values of the weights, $\boldsymbol{\beta}$, and the **hyperparameter** $\gamma \in \mathbb{R}_+$ sets the strength of the penalty. Even in the linearly separable regime, the maximizer is finite.

Now, the gradient and Hessian are,

$$\begin{aligned}\nabla \mathcal{L}(\boldsymbol{\beta}) &= -\frac{1}{n} \sum_{i=1}^n (y_i - \sigma(\mathbf{x}_i^\top \boldsymbol{\beta})) \mathbf{x}_i + \gamma \boldsymbol{\beta} \\ \nabla^2 \mathcal{L}(\boldsymbol{\beta}) &= \frac{1}{n} \sum_{i=1}^n w_i \mathbf{x}_i \mathbf{x}_i^\top + \gamma \mathbf{I}.\end{aligned}$$

Choosing the Hyperparameter

There are many ways to select a value of γ .

One approach is to *not* choose a single value and instead try to compute the **regularization path**; i.e., the solution $\hat{\beta}(\gamma)$ for a range of $\gamma \in [0, \infty)$.

Another is to hold out a fraction of data and use cross-validation to select the hyperparameter setting that yields the best performance on the held-out data.

Bayesian Perspective

From a Bayesian perspective, we can think of the regularizer as a **prior log probability**.

Here, the regularizer corresponds to a spherical Gaussian prior,

$$\boldsymbol{\beta} \sim \mathcal{N}(\mathbf{0}, (\gamma n)^{-1} \mathbf{I}),$$

where **precision (inverse covariance)** γn .

Minimizing the objective above corresponds to doing **maximum a posteriori (MAP)** estimation in the Bayesian model.

We'll talk more about Bayesian methods in the coming weeks.

Newton's Method

Gradient descent leverages the gradient at β to determine the update. Newton's method uses the Hessian to inform the update as well, and in doing so it can achieve considerably faster convergence.

The second order Taylor approximation of \mathcal{L} around β is

$$\mathcal{L}(\beta') \approx \mathcal{L}(\beta) + \nabla \mathcal{L}(\beta)^\top (\beta' - \beta) + \frac{1}{2} (\beta' - \beta)^\top \nabla^2 \mathcal{L}(\beta) (\beta' - \beta) \triangleq \widetilde{\mathcal{L}}(\beta')$$

The stationary point of $\widetilde{\mathcal{L}}(\beta')$ is at

$$\begin{aligned} \nabla \widetilde{\mathcal{L}}(\beta') &= \nabla \mathcal{L}(\beta) + \nabla^2 \mathcal{L}(\beta) (\beta' - \beta) = 0 \\ \implies \beta' &= \beta - [\nabla^2 \mathcal{L}(\beta)]^{-1} \nabla \mathcal{L}(\beta), \end{aligned}$$

assuming the Hessian is invertible. When $\nabla^2 \mathcal{L}(\beta) \succ 0$ – i.e., when the Hessian is positive definite – the inverse Hessian exists and the stationary point is a minimizer.

Newton's Method

Vanilla Newton's method applies this update repeatedly until convergence, forming a quadratic approximation and then minimizing it.

Damped Newton's method adds a step size $\alpha_t < 1$ to improve stability,

$$\beta_{t+1} \leftarrow \beta_t - \alpha_t [\nabla^2 \mathcal{L}(\beta_t)]^{-1} \nabla \mathcal{L}(\beta_t),$$

The step size can be chosen by backtracking line search.

Question: Compare the Newton update to that of gradient descent. How do they differ?

Iteratively Reweighted Least Squares

Newton's method can be viewed as iteratively solving a weighted least squares problem.

Write the gradient and Hessian in matrix form,

$$\begin{aligned}\nabla \mathcal{L}(\boldsymbol{\beta}_t) &= -\frac{1}{n} \mathbf{X}^\top (\mathbf{y} - \hat{\mathbf{y}}_t) \\ \nabla^2 \mathcal{L}(\boldsymbol{\beta}_t) &= \frac{1}{n} \mathbf{X}^\top \mathbf{W}_t \mathbf{X}\end{aligned}$$

where

- ▶ $\mathbf{X} \in \mathbb{R}^{n \times p}$ is the design matrix with rows \mathbf{x}_i^\top
- ▶ $\mathbf{y} \in \{0, 1\}^n$ is the vector of binary responses
- ▶ $\hat{\mathbf{y}}_t = \sigma(\mathbf{X}\boldsymbol{\beta}_t) \in [0, 1]^n$ is the vector of predicted response means
- ▶ $\mathbf{W}_t = \text{diag}([w_{t,1}, \dots, w_{t,n}])$ with $w_{t,i} = \sigma(\mathbf{x}_i^\top \boldsymbol{\beta}_t)(1 - \sigma(\mathbf{x}_i^\top \boldsymbol{\beta}_t))$ is the diagonal weight matrix.

Iteratively Reweighted Least Squares

Then the (undamped) Newton update is,

$$\begin{aligned}\beta_{t+1} &= \beta_t - [\nabla^2 \mathcal{L}(\beta_t)]^{-1} \nabla \mathcal{L}(\beta_t) \\ &= \beta_t + [\mathbf{X}^\top \mathbf{W}_t \mathbf{X}]^{-1} \mathbf{X}^\top (\mathbf{y} - \hat{\mathbf{y}}_t) \\ &= [\mathbf{X}^\top \mathbf{W}_t \mathbf{X}]^{-1} \mathbf{X}^\top \mathbf{W}_t \mathbf{X} \beta_t + [\mathbf{X}^\top \mathbf{W}_t \mathbf{X}]^{-1} \mathbf{X}^\top (\mathbf{y} - \hat{\mathbf{y}}_t) \\ &= [\mathbf{X}^\top \mathbf{W}_t \mathbf{X}]^{-1} \mathbf{X}^\top \mathbf{W}_t \mathbf{z}_t\end{aligned}$$

where

$$\mathbf{z}_t = \mathbf{X} \beta_t + \mathbf{W}_t^{-1} (\mathbf{y} - \hat{\mathbf{y}}_t).$$

This is a **weighted least squares** problem with **weights** $w_{t,j}$ and **adjusted (or working) responses** $z_{t,j}$.

Iteratively Reweighted Least Squares

How can we interpret the working responses? We can view them as the real responses mapped through a Taylor approximation of the link (inverse mean) function,

$$\begin{aligned}\sigma^{-1}(y_i) &\approx \sigma^{-1}(\hat{y}_{t,i}) + (y_i - \hat{y}_{t,i})[\sigma^{-1}]'(\hat{y}_{t,i}) \\ &= \mathbf{x}_i^\top \boldsymbol{\beta}_t + \frac{(y_i - \hat{y}_{t,i})}{w_{t,i}} \\ &\triangleq z_{t,i}.\end{aligned}$$

Asymptotic Covariance of MLE

What other insight can we glean from the Hessian? Recall our discussion of the asymptotic normality of the MLE from Lecture 1. For iid observations $Y_i \stackrel{\text{iid}}{\sim} p(\cdot; \boldsymbol{\theta})$ for $i = 1, \dots, n$, the asymptotic covariance is $\mathcal{I}(\boldsymbol{\theta})^{-1}/n$, where

$$\mathcal{I}(\boldsymbol{\theta}) = -\mathbb{E}[\nabla_{\boldsymbol{\theta}}^2 \log p(Y; \boldsymbol{\theta})]$$

is the Fisher information matrix.

For logistic regression, we have n independent but *not* identically distributed observations. In this case, the asymptotic covariance follows the same form. It is given by the inverse of the Fisher information matrix; i.e., the inverse of the negative expected Hessian of the log likelihood,

$$\mathcal{I}(\boldsymbol{\beta}) = -\sum_{i=1}^N \mathbb{E}[\nabla_{\boldsymbol{\beta}}^2 \log p(Y_i | \mathbf{X}_i = \mathbf{x}_i; \boldsymbol{\beta})].$$

(Note that this includes the iid formula as a special case.)

Asymptotic Covariance of MLE

Substituting the form of the Hessian from above,

$$\begin{aligned}\mathcal{J}(\boldsymbol{\beta}) &= \sum_{i=1}^n \mathbb{E}[\text{Var}[Y_i | \mathbf{X} = \mathbf{x}_i] \mathbf{x}_i \mathbf{x}_i^\top] \\ &= \sum_{i=1}^n w_i \mathbf{x}_i \mathbf{x}_i^\top\end{aligned}$$

where $w_i = \sigma(\boldsymbol{\beta}^\top \mathbf{x}_i)(1 - \sigma(\boldsymbol{\beta}^\top \mathbf{x}_i))$. Finally, we evaluate the Fisher information the MLE to obtain the asymptotic covariance estimate,

$$\widehat{\text{Cov}}(\hat{\boldsymbol{\beta}}) = [\mathcal{J}(\hat{\boldsymbol{\beta}})]^{-1}.$$

Like before, we can use the asymptotic covariance estimate to derive Wald confidence intervals for the parameters and perform hypothesis tests.

Asymptotic Covariance of MLE

Caution: Remember, Wald confidence intervals are only as good as the asymptotic normality assumption. When the likelihood is not well approximated by a quadratic, the covariance estimate will be poor, and the confidence intervals will be invalid. When might the Gaussian approximation not hold? :::

Converge Rate of Gradient Descent

To determine when and at what rate gradient descent converges, we need to know more about the eigenvalues of the Hessian.

If we can bound the maximum eigenvalue of the Hessian by L , then we can obtain a quadratic upper bound on the negative log likelihood,

$$\mathcal{L}(\beta') \leq \mathcal{L}(\beta) + \nabla \mathcal{L}(\beta)^\top (\beta' - \beta) + \frac{L}{2} (\beta' - \beta)^\top (\beta' - \beta).$$

That means the negative log likelihood is an L -smooth function.

Converge Rate of Gradient Descent

Example: Bounded Covariates

If the covariates have bounded norm, $\|\mathbf{x}_i\|_2 \leq B$, then we can bound the maximum eigenvalue of the Hessian by,

$$\begin{aligned}\lambda_{\max} &= \max_{\mathbf{u} \in \mathbb{S}_{p-1}} \mathbf{u}^\top \nabla^2 \mathcal{L}(\boldsymbol{\beta}) \mathbf{u} \\ &= \max_{\mathbf{u} \in \mathbb{S}_{p-1}} \frac{1}{n} \sum_{i=1}^n w_i \mathbf{u}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{u} \\ &\leq \frac{B^2}{4}\end{aligned}$$

since the weights are the conditional variances of Bernoulli random variables, which are at most $\frac{1}{4}$, and since $\mathbf{u}^\top \mathbf{x}_i \leq B$ for all unit vectors $\mathbf{u} \in \mathbb{S}_{p-1}$ (the unit sphere embedded in \mathbb{R}^p). This isn't meant to be a tight upper bound.

Converge Rate of Gradient Descent

If we run gradient descent with a constant step size of $\alpha = 1/L$, then the algorithm converges at a rate of $1/t$, which means that after t iterations

$$\mathcal{L}(\boldsymbol{\beta}^*) - \mathcal{L}(\boldsymbol{\beta}_t) \leq \frac{L}{t} \|\boldsymbol{\beta}^* - \boldsymbol{\beta}_0\|_2^2,$$

where $\boldsymbol{\beta}_0$ is the initial setting of the parameters and $\boldsymbol{\beta}^*$ is the global optimum.

Put differently, if we want a gap of at most epsilon, we need to run $t \sim \mathcal{O}(1/\epsilon)$ iterations of gradient descent. Put differently, if we want to reduce ϵ by a factor of 100, we need to run around 100 times as many iterations. This is called a **sub-linear convergence** rate.

Converge Rate of Gradient Descent with Regularization

With regularization, we can also **lower bound** the objective by a quadratic function,

$$\mathcal{L}(\beta') \geq \mathcal{L}(\beta) + \nabla \mathcal{L}(\beta)^\top (\beta' - \beta) + \frac{\mu}{2} (\beta' - \beta)^\top (\beta' - \beta)$$

for $\mu > 0$, which means the objective is μ -**strongly convex**.

For twice differentiable objectives, the minimum eigenvalue of Hessian provides such a lower bound, $\mu = \lambda_{\min}$. In particular, we know that minimum eigenvalue of $\nabla^2 \mathcal{L}(\beta)$ is at least γ . (This bound is achieved when the data are linearly separable, the model is saturated, and the conditional variances are all zero.)

For a L -smooth and μ -strongly convex function with stepsize $\alpha = 1/L$, gradient descent has the following convergence guarantee,

$$\mathcal{L}(\beta_{t+1}) - \mathcal{L}(\beta^*) \leq \left(1 - \frac{\mu}{L}\right) (\mathcal{L}(\beta_t) - \mathcal{L}(\beta^*)).$$

Converge Rate of Gradient Descent with Regularization

Applying this bound recursively yields that,

$$\mathcal{L}(\beta_t) - \mathcal{L}(\beta^*) \leq \left(1 - \frac{\mu}{L}\right)^t (\mathcal{L}(\beta_0) - \mathcal{L}(\beta^*)).$$

If we want to find the number of iterations t to bound the gap by at most ϵ , we need to solve for t in

$$\left(1 - \frac{\mu}{L}\right)^t (\mathcal{L}(\beta_0) - \mathcal{L}(\beta^*)) \leq \epsilon.$$

This inequality is equivalent to taking the log of both sides

$$\log(\mathcal{L}(\beta_0) - \mathcal{L}(\beta^*)) + t \log\left(1 - \frac{\mu}{L}\right) \leq \log \epsilon.$$

Converge Rate of Gradient Descent with Regularization

We can further upper bound the LHS by using the inequality $\log(1 - x) \leq -x$ to get

$$\log(\mathcal{L}(\beta_0) - \mathcal{L}(\beta^*)) - \frac{\mu t}{L} \leq \log \epsilon. \quad (1)$$

So, we need to run $t \geq \frac{L}{\mu} \log \frac{\mathcal{L}(\beta_0) - \mathcal{L}(\beta^*)}{\epsilon} \sim \log \frac{1}{\epsilon}$ iterations of gradient descent. If we want to reduce ϵ by a factor of 100, we only need to run around $\log 100$ times as many iterations. This is called **linear convergence**.

Converge Rate of Newton's Method

Under certain conditions – if the objective is strongly convex, the Hessian is Lipschitz continuous, and we start near enough to the global optimum – Newton's method achieves **second order convergence**, meaning

$$\|\beta_{t+1} - \beta^*\|_2 \leq (c\|\beta_t - \beta^*\|_2)^2$$

for some positive constant c , provided we start with β_0 close enough to β^* . Applying this bound recursively yields that,

$$\|\beta_t - \beta^*\|_2 \leq (c\|\beta_0 - \beta^*\|_2)^{2^t}.$$

Put differently, if we start with $\|\beta_0 - \beta^*\| < c^{-1}$, then we need $t \sim \mathcal{O}(\log \log \frac{1}{\epsilon})$ iterations to obtain an error of ϵ . Since the double log grows incredibly slowly, this statement says that we effectively need a constant number of iterations for Newton's method to converge in this regime.

For more information on convergence rates of gradient descent and Newton's method, see [boyd2004convex](#), ch. 9.

Conclusion

One, two, many... with logistic regression, we can begin modeling relationships between binary response variables and (possibly arbitrary-valued) covariates. Next, we'll see how to expand from logistic regression to generalized linear models for exponential family responses.