

Quiz 1 - DSCI 562

Review of Regression I

Ordinary Least Squares (OLS) Regression

- Response of continuous nature (e.g., "Salary")
- Response is subject to regressors (or explanatory variables/ features/ independent variables)
- More than 1 regressor \Rightarrow multiple linear regression

$$Y_i = \beta_0 + \beta_1 g_1(x_{i1}) + \beta_2 g_2(x_{i2}) + \cdots + \beta_p g_p(x_{ip}) + \epsilon_i$$

ϵ_i is the error term

- 1. Linearity: the relationship between the response and functions of the regressors is linear
- 2. Errors are independent of each other and are normally distributed with mean 0 and variance σ^2

Hence, each Y_i is assumed to be independent and normally distributed.

- To fit we'll need $k+2$ parameters: $\beta_0, \beta_1, \dots, \beta_k, \sigma^2$
- Minimize the sum of squared errors (SSE) OR maximize the likelihood of the observed data

- Maximum Likelihood Estimation (MLE): find the parameters that maximize the likelihood of the observed data
 - Likelihood: the probability of observing the data given the parameters
 - Log-Likelihood: the log of the likelihood

- Do a t-test on the parameters to see if they are statistically significant

Limitations of OLS

- Ols allows responses to take any real number.
- Example of non-continuous responses:
 - Non-negative values
 - Binary values (success/failure)
 - Count data

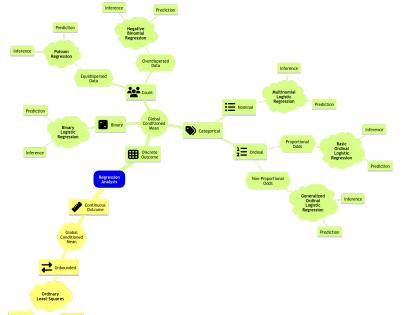
Link Function

- Recall OLS models a continuous response via its conditional mean
- $\mu = E(Y|X_i) = \beta_0 + \beta_1 g_1(x_{i1}) + \beta_2 g_2(x_{i2}) + \cdots + \beta_p g_p(x_{ip})$
- BUT this is not suitable for non-continuous responses (e.g. binary, count, non-negative).
- Solution: use a link function $h(\mu_i)$ to map the conditional mean to the real line
- Link function: relate the systematic component, η_i , with the response's mean
- $h(\mu_i) = \eta_i = \beta_0 + \beta_1 g_1(x_{i1}) + \beta_2 g_2(x_{i2}) + \cdots + \beta_p g_p(x_{ip})$
- Monotonic: allows for a one-to-one mapping between the mean of the response variable and the linear predictor
- Differentiable: to allow for maximum likelihood estimation (MLE), used to obtain $\hat{\beta}$

Generalized Linear Models (GLM)

- Generalized Linear Models (GLM): a generalization of OLS regression that allows for non-continuous responses

GLM = link function + error distribution



Poisson Regression

- Poisson regression: a GLM for count data (Equidispersed)
- Equidispersed: the variance of the response is equal to its mean (i.e. $Var(Y_i) = E(Y_i) = \lambda_i$)
- It assumes a random sample of n count observations Y_i
- Independent
- Not identically distributed: Each Y_i has its own mean $E(Y_i) = \lambda_i > 0$ and variance $Var(Y_i) = \lambda_i > 0$

$$Y_i \sim Poisson(\lambda_i)$$

- λ_i is the risk of event occurrence in a given timeframe or area (definition of Poisson distribution)

Link Function for Poisson Regression

- Log Link function: the log of the mean of the response variable is linearly related to the regressors

$$h(\mu_i) = \log(\mu_i) = \eta_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \cdots + \beta_p X_{ip}$$

Hence,

$$\lambda_i = e^{\eta_i} = e^{\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \cdots + \beta_p X_{ip}}$$

- This is good since λ_i (mean count) is always positive

Poisson Regression in R

```
glm(Y ~ X, family = poisson, data = dataset)

# view each regression coefficient
tidy(glm_model)
tidy(glm_model, confint = TRUE) # for 95% confidence interval

# view model summary
glance(glm_model)
```

Interpretation of Coeffs of Poisson Regression

e.g. $\beta_1 = 0.5$

- β_1 is the expected change in the log of the mean count for a one-unit increase in X_1 holding all other variables constant
- one-unit increase in X_1 will increase the mean count by $e^{0.5} = 1.65$ times.

Inference of Poisson Regression

- To determine the significance of the parameters $\beta_1, \beta_2, \dots, \beta_k$, we can do a Wald statistic

$$z_j = \frac{\beta_j}{SE(\beta_j)}$$

To test the hypothesis:

- $H_0: \beta_j = 0$
- $H_1: \beta_j \neq 0$

Negative Binomial Regression

- Negative Binomial Regression: a GLM for count data (Overdispersed)

- Oversdispersion: the variance of the response is greater than its mean (i.e. $Var(Y_i) > E(Y_i) = \lambda_i$)

Check for Overdispersion

dispersiontest(glm_model)

- If $p\text{-value} < 0.05$, then there is overdispersion (reject null hypothesis)

If use Poisson regression on overdispersed data, then the standard errors will be underestimated \Rightarrow Type I error (false positive) increases

Recall PMF of Negative Binomial Distribution:

$$P(Y_i | m, p_i) = \binom{m+i-1}{i} p_i^m (1-p_i)^i$$

y_i is the number of failures before experiencing m successes where probability of success is p_i

$$E(Y_i) = \frac{m(1-p_i)}{p_i}$$

$$\text{Var}(Y_i) = \frac{m(1-p_i)}{p_i^2}$$

Rearranging the above equations, we get:

$$E(Y_i) = \lambda_i$$

$$\text{Var}(Y_i) = \lambda_i(1 + \frac{\lambda_i}{m})$$

Interesting information:

$$X \sim Poisson(\lambda) = \lim_{m \rightarrow \infty} \text{NegativeBinomial}(m, p)$$

Negative Binomial Regression in R

glm.nb(Y ~ X, data = dataset)

- RE: needs to select model with fewer regressors than AIC.
- smaller AIC is better (means better fit)
- Can get from `glance` function

Multinomial Logistic Regression

- Is a MLE-based GLM when the response is categorical and nominal.
 - Nominal: unordered categories
 - e.g. red, green, blue
 - Ordinal: ordered categories
 - e.g. low, medium, high
- Similar to binomial logistic regression, but with more than 2 categories.
- Levels: levels of the response variable
 - need more than 1 logit function to model the probabilities of each category.
 - One category is the baseline category, the other categories are compared to the baseline category.

$$\eta_1^{(model1, model2)} = \log \left[\frac{P(Y_i = 1 | X_{i1}, X_{i2}, X_{i3})}{P(Y_i = 0 | X_{i1}, X_{i2}, X_{i3})} \right] = \beta_0^{(model1, model2)} + \beta_1^{(model1, model2)} X_{i1} + \beta_2^{(model1, model2)} X_{i2} + \beta_3^{(model1, model2)} X_{i3}$$

$$\eta_2^{(model1, model2)} = \log \left[\frac{P(Y_i = 2 | X_{i1}, X_{i2}, X_{i3})}{P(Y_i = 0 | X_{i1}, X_{i2}, X_{i3})} \right] = \beta_0^{(model1, model2)} + \beta_1^{(model1, model2)} X_{i1} + \beta_2^{(model1, model2)} X_{i2} + \beta_3^{(model1, model2)} X_{i3}$$

With some algebra, we can get the following (for M categories):

$$p_{i,1}^{(model1)} = \frac{1}{1 + \sum_{j=2}^M e^{\eta_j^{(model1, model2)}}}$$

$$p_{i,2}^{(model1)} = \frac{e^{\eta_2^{(model1, model2)}}}{1 + \sum_{j=2}^M e^{\eta_j^{(model1, model2)}}}$$

All probabilities sum to 1.

Nuances: Baseline Category

- The baseline level is the level that is not included in the model
- Can find using `levels` function, the first levels is the baseline level.

levels(data\$response) # to check levels

to change levels:

- `base = "new_level_W"`
- `"W" = "new_level_W"`
- `"T" = "new_level_L"`

Estimation of MLR

model = multinom(response ~ regressor_1 + regressor_2 + regressor_3, data = data)

to get test statistics

mlr_output = mlr_test(model)

confint = mlr_output\$confint # to get confidence intervals (default is 95%)

exponentiate = TRUE # to get odds ratios

default result is log odds ratios

can filter p-values

mlr_output\$p.value # filter(p.value < 0.05)

predict

predict(model, newdata = data, type = "probs")

sum of all probabilities is 1

Inference of MLR

- Check if regressor is significant using `waldTest`.

glm.nb(Y ~ X, data = dataset)

Check if regressor is significant using `WaldTest`.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

.

- j is the level of the response variable

$$\sum_{j=1}^J p_{ij} = 1$$

Estimation of OLR

```
• use MASS::polr function
ordinal_model <- polr(
  formula = response ~ regressor_1 + regressor_2,
  data = data,
  Hess = TRUE # Hessian matrix of log-likelihood
)
```

Inference of OLR

- Similar to LRL using `wald.test`

```
ckind
tidy(ordinal_model),
p.value = pnorm(absolute(tidy(ordinal_model)$statistic),
  1) < 2
)
# confidence intervals
confint(ordinal_model) # default is 95%
```

Coefficient Interpretation of OLR

- e.g. $\beta_1 = 0.6$
 - For a one unit increase in X_1 , the odds of being in a higher category is $e^{0.6} = 1.82$ times the odds of being in a lower category, holding all other variables constant.

Predictions

```
predict(ordinal_model, newdata = data, type = "prob")
# returns probabilities for each level
```

- To get the corresponding predicted cumulative odds for a new observation, use `VGAM::vglm` function

```
otr <- vglm(
  response ~ regressor_1 + regressor_2,
  propodds, # for proportional odds model
  data,
)
# can also predict using this model, same as code block above
predict(otr, newdata = data, type = "response")
predict(otr, newdata = data, type = "link") |>
  exp() # to get odds instead of log odds

• Interpret the predicted cumulative odds as:
  • e.g.  $logit(\text{Pr}(Y_j \geq j)) = 2.68$ 

- A student with  $X_1$  is 2.68 times more likely to be in  $j$  or higher category than in category  $j - 1$  or lower, holding all other variables constant.


  •  $exp(2.68) = 13.63$ 

- A student with  $X_1$  is 3.03 (0.33) times more likely to be in  $j$  category or lower than in category  $j$  or higher, holding all other variables constant.

```

Non-proportional Odds

- If the proportional odds assumption is not met, we can use the **generalized ordinal logistic regression model**.
 - If the proportional odds assumption is not met, we can use the **generalized ordinal logistic regression model**.
 - Basically all β 's are allowed to vary across the different levels of the response variable.

Brant(ordinal_model)

- If the proportional odds assumption is not met, we can use the **generalized ordinal logistic regression model**.

```
• Test for proportional odds assumption using the Brant-Null test
  •  $H_0$ : Our OLR model globally fulfills the proportional odds assumption.
  •  $H_A$ : Our OLR model does not globally fulfill the proportional odds assumption.
```

Linear Fixed Effects Model

- Linear Fixed Effects Model (LFE) is a generalization of the linear regression model
- Fixed Effects: the parameters of the model
 - constant for all observations

Limitations

- Data Hierarchy: the data is organized in a hierarchy
 - Can be due to **sampling levels**
 - e.g. investments in different firms, students in different schools (sampling schemes may be different in different schools)
 - Might have some correlation between datapoints in firms/ schools
 - Violates the independence assumption (I.i.d. observations)

Example: Investments in different firms

- Goal: assessing the association of gross investment with market_value and capital in the population of American firms.
- Date 11 firm, 20 observations per firm
 - 2 hierarchical levels: firm and observation

1. Trial 1: ignore hierarchy

```
ordinary_model <- lm(
  formula = investment ~ market_value + capital,
  data = Grunfeld)
```

2. Trial 2: Different intercepts for different firms

```
model_waring_intercept <- lm(
  # i.e. so that baseline is not included as first intercept.
  formula = investment ~ market_value + capital + firm = 1,
  data = Grunfeld)
```

3. Trial 3: OLS regression for each firm

- This does NOT solve our goal
- We want to find out among all firms, not one specific firm.

```
model_by_firm <- lm(
  investment ~ market_value + firm + capital + firm,
  data = Grunfeld)
```

Linear Mixed Effects Model

- Fundamental idea:
 - data subsets of elements share a **common structure**
 - i.e. all n rows of training data are not independent

mixed effect = fixed effect + random effect

$$\beta_{0j} = \beta_0 + b_{0j}$$

- β_{0j} mixed effect: the intercept for the j th school/firm

- β_0 fixed effect: the average intercept

- b_{0j} random effect: the deviation of the j th school/firm from the average intercept

- $b_{0j} \sim N(0, \sigma_b^2)$

- Independent of the error term ϵ

- Variance of the i th observation:
 - $\sigma_b^2 + \sigma_e^2$

Full Equation for LME

$$y_{ij} = \beta_{0j} + \beta_{1j}x_{1ij} + \beta_{2j}x_{2ij} + \epsilon_{ij} = (\beta_0 + b_{0j}) + (\beta_1 + b_{1j})x_{1ij} + (\beta_2 + b_{2j})x_{2ij} + \epsilon_{ij}$$

Note: $(\beta_{0j}, b_{0j}, b_{1j}, b_{2j}) \sim N(0, D)$

- D: vector of zero, e.g. $(0, 0, 0)^T$
- D: generic covariance matrix

$$D = \begin{bmatrix} \sigma_0^2 & \sigma_{01} & \sigma_{02} \\ \sigma_{10} & \sigma_1^2 & \sigma_{12} \\ \sigma_{20} & \sigma_{21} & \sigma_2^2 \end{bmatrix} = \begin{bmatrix} \sigma_0^2 & \rho_{01}\sigma_0\sigma_1 & \rho_{02}\sigma_0\sigma_2 \\ \rho_{10}\sigma_0\sigma_1 & \sigma_1^2 & \rho_{12}\sigma_1\sigma_2 \\ \rho_{20}\sigma_0\sigma_2 & \rho_{21}\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}$$

- ρ_{uv} : pearson correlation between u th and v th random effects

Model Fitting of LME

- use the `lme` function from the `lme4` package

```
mixed_intercept_model <- lmer(
  response ~ regressor_1 + regressor_2 +
  (1 | school), # random intercept by firm
  data
)

full_model <- lmer(
  response ~ regressor_1 + regressor_2 +
  (regressor_1 + regressor_2 | school),
  # random intercept and slope by firm
  data
)
```

- Equation for mixed intercept model:

$$y_{ij} = (\beta_0 + b_{0j}) + \beta_1x_{1ij} + \beta_2x_{2ij} + \epsilon_{ij}$$

- Equation for full model:

$$y_{ij} = (\beta_0 + b_{0j}) + (\beta_1 + b_{1j})x_{1ij} + (\beta_2 + b_{2j})x_{2ij} + \epsilon_{ij}$$

Inference of LME

- Cannot do inference using normal t-test

```
summary(fixed_intercept_model)
summary(full_model)

# obtain coefficients
coeff(fixed_intercept_model)$tstat
coeff(full_model)$tstat
```

Prediction with LME

- Predict on existing group

- Predict on new group

```
predict(full_model,
  newdata = tibble(school = "new_school", regressor_1 = 1, regressor_2 = 2))
```

Brant(ordinal_model)

- If the proportional odds assumption is not met, we can use the **generalized ordinal logistic regression model**.

```
• If the proportional odds assumption is not met, we can use the generalized ordinal logistic regression model.
  • Basically all  $\beta$ 's are allowed to vary across the different levels of the response variable.
```

Linear Fixed Effects Model

- Linear Fixed Effects Model (LFE) is a generalization of the linear regression model
- Fixed Effects: the parameters of the model
 - constant for all observations

Limitations

- Data Hierarchy: the data is organized in a hierarchy
 - Can be due to **sampling levels**
 - e.g. investments in different firms, students in different schools (sampling schemes may be different in different schools)
 - Might have some correlation between datapoints in firms/ schools
 - Violates the independence assumption (I.i.d. observations)

Example: Investments in different firms

- Goal: assessing the association of gross investment with market_value and capital in the population of American firms.
- Date 11 firm, 20 observations per firm
 - 2 hierarchical levels: firm and observation

1. Trial 1: ignore hierarchy

```
ordinary_model <- lm(
  formula = investment ~ market_value + capital,
  data = Grunfeld)
```

2. Trial 2: Different intercepts for different firms

```
model_waring_intercept <- lm(
  # i.e. so that baseline is not included as first intercept.
  formula = investment ~ market_value + capital + firm = 1,
  data = Grunfeld)
```

3. Trial 3: OLS regression for each firm

- This does NOT solve our goal
- We want to find out among all firms, not one specific firm.

```
model_by_firm <- lm(
  investment ~ market_value + firm + capital + firm,
  data = Grunfeld)
```

Linear Mixed Effects Model

- Fundamental idea:
 - data subsets of elements share a **common structure**
 - i.e. all n rows of training data are not independent

mixed effect = fixed effect + random effect

$$\beta_{0j} = \beta_0 + b_{0j}$$

- β_{0j} mixed effect: the intercept for the j th school/firm

- β_0 fixed effect: the average intercept

- b_{0j} random effect: the deviation of the j th school/firm from the average intercept

- $b_{0j} \sim N(0, \sigma_b^2)$

- Independent of the error term ϵ

- Variance of the i th observation:
 - $\sigma_b^2 + \sigma_e^2$