

Generalized Multilevel Models for Non-Normal Longitudinal Data

- Topics:
 - **3 parts of a generalized (multilevel) model**
 - Models for binary outcomes
 - Models for categorical outcomes
 - Complications for generalized multilevel models
 - A brief tour of other generalized models:
 - Models for count outcomes
 - Models for not-normal but continuous outcomes

Dimensions for Organizing Models

- Outcome type: General (normal) vs. Generalized (not normal)
- Dimensions of sampling: One (so one variance term per outcome) vs. **Multiple** (so multiple variance terms per outcome) → **OUR WORLD**
- **General Linear Models**: conditionally normal outcome distribution, **fixed effects** (identity link; only one dimension of sampling)
- **Generalized Linear Models**: **any conditional outcome distribution**, **fixed** effects through **link functions**, no random effects (one dimension)
- **General Linear Mixed Models**: conditionally normal outcome distribution, **fixed and random effects** (identity link, but multiple sampling dimensions)
- **Generalized Linear Mixed Models**: **any conditional outcome distribution**, **fixed and random effects** through **link functions** (multiple dimensions)
- “Linear” means the fixed effects predict the *link-transformed* DV in a linear combination of (effect*predictor) + (effect*predictor)...

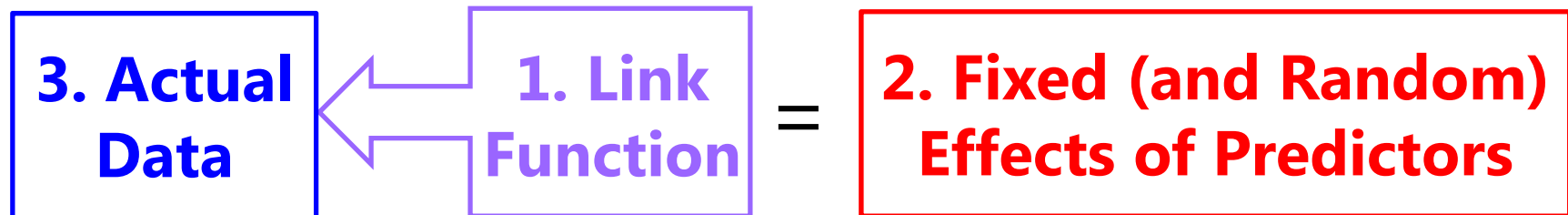
Note: Least Squares is only for GLM

Generalized Linear Models

- **Generalized linear models:** link-transformed Y is predicted instead of actual Y ; ML estimator uses not-normal distributions
 - **Single-level models** → residuals follow some not-normal distribution
 - **Multilevel models** → level-1 residuals follow some not-normal distribution, but level-2 random effects are almost always still MVN
- Many kinds of non-normally distributed outcomes have some kind of generalized linear model to go with them via ML:
 - Binary (dichotomous)
 - Unordered categorical (nominal)
 - Ordered categorical (ordinal)
 - Counts (discrete, positive values)
 - Censored (piled up and cut off at one end)
 - Zero-inflated (pile of 0's, then some distribution after)
 - Continuous but skewed data (long tail)

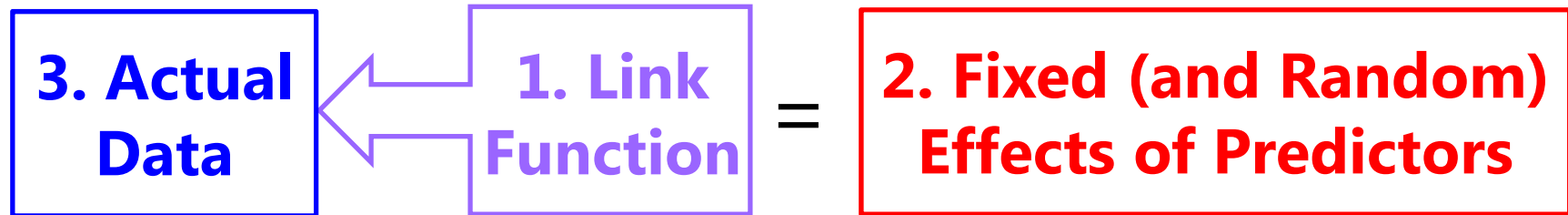
} These two are often called
"multinomial" inconsistently

3 Parts of Generalized Multilevel Models



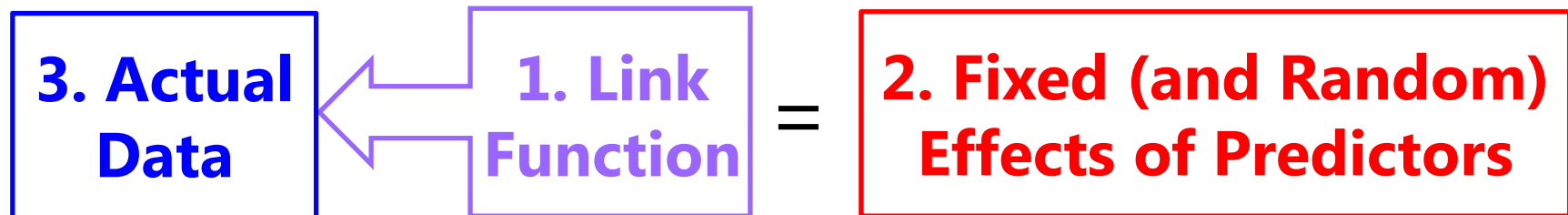
1. Link Function (different from general): How a non-normal outcome is transformed into an **unbounded** outcome that the model fixed and random effects can predict linearly
 - Transformed outcome is predicted directly, then converted back into Y
 - This way the predicted outcomes will stay within the sample space (boundaries) of the observed data (e.g., 0/1 for binary outcomes—the model should not predict -1 or 2 , so linear slopes need to shut off)
 - Written as $g(\cdot)$ for link and $g^{-1}(\cdot)$ for inverse link (to go back to data)
 - For outcomes with residuals that are already normal, general linear models are just a special case with an “identity” link function ($Y * 1$)
 - So general linear models are a special case of *generalized* linear models, and general linear mixed models are a special case of *generalized* linear mixed models

3 Parts of Generalized Multilevel Models



2. **Linear Predictor** (same as in general): How the model predictors linearly relate to the transformed outcome
- This works the same as usual, except the linear predictor model **directly predicts the link-transformed outcome**, which then gets converted back into the original outcome
 - That way we can still use the familiar “one-unit change” language to describe the effects of model predictors
 - You can think of this as “model for the means” still, but it also includes the level-2 random effects for dependency of level-1 observations
 - Fixed effects are no longer determined: they now have to be found through the ML algorithm, the same as the variance parameters

3 Parts of Generalized Multilevel Models



3. **Model for Level-1 Residuals** (different than general): how level-1 residuals should be distributed given the sample space of the actual outcome
- Many alternative distributions that map onto what the distribution of **residuals** could possibly look like (and kept within sample space)
 - **Why?** To get the most correct **standard errors** for fixed effects
 - You can think of this as “model for the variance” still, but not all distributions have estimated residual variance
- Let's start with models for **binary data** to illustrate these 3 parts...

Generalized Multilevel Models for Non-Normal Longitudinal Data

- Topics:
 - 3 parts of a generalized (multilevel) model
 - **Models for binary outcomes**
 - Models for categorical outcomes
 - Complications for generalized multilevel models
 - A brief tour of other generalized models:
 - Models for count outcomes
 - Models for not-normal but continuous outcomes

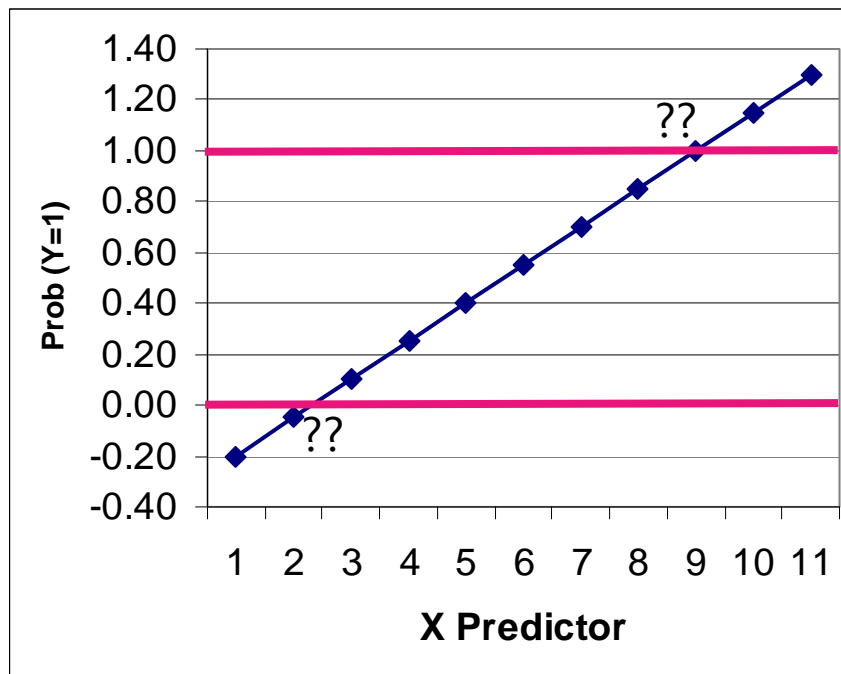
Normal GLM for Binary Outcomes?

- Let's say we have a single binary (0 or 1) outcome...
(*concepts for longitudinal data will proceed similarly*)
 - Expected mean is proportion of people who have a 1, so the **probability of having a 1** is what we're trying to predict for each person, given the predictor values: $P(y_i = 1)$
 - General linear model: $P(y_i = 1) = \beta_0 + \beta_1 X_i + \beta_2 Z_i + e_i$
 - β_0 = expected probability when all predictors are 0
 - β 's = expected change in $P(y_i = 1)$ for a one-unit Δ in predictor
 - e_i = difference between observed and predicted binary values
 - Model becomes $y_i = (\text{predicted probability of 1}) + e_i$
 - **What could possibly go wrong?**

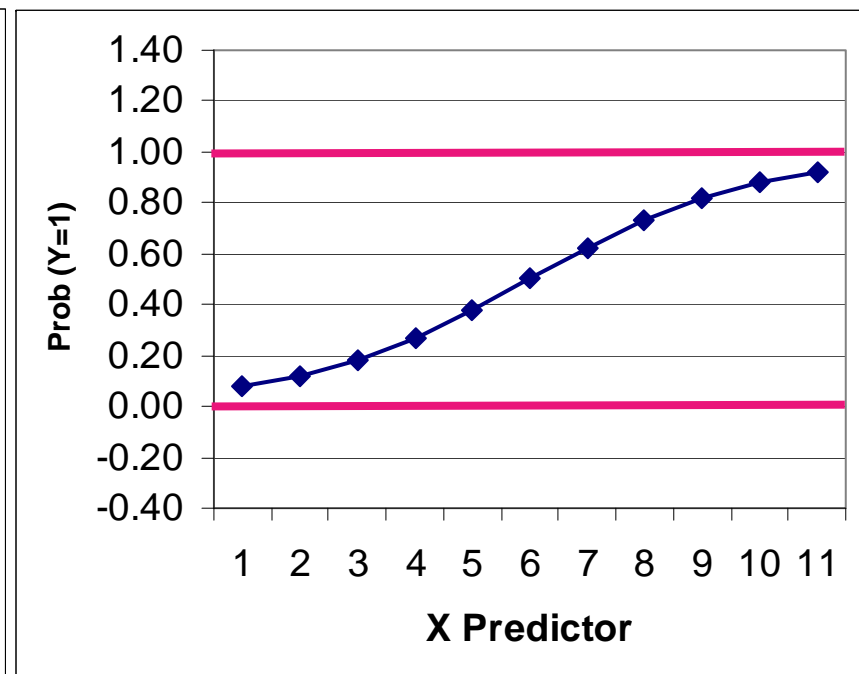
Normal GLM for Binary Outcomes?

- Problem #1: A **linear** relationship between X and Y???
- Probability of a 1 is bounded between 0 and 1, but predicted probabilities from a linear model aren't bounded
- Linear relationship needs to shut off → made nonlinear

We have this...



But we need this...



Generalized Models for Binary Outcomes

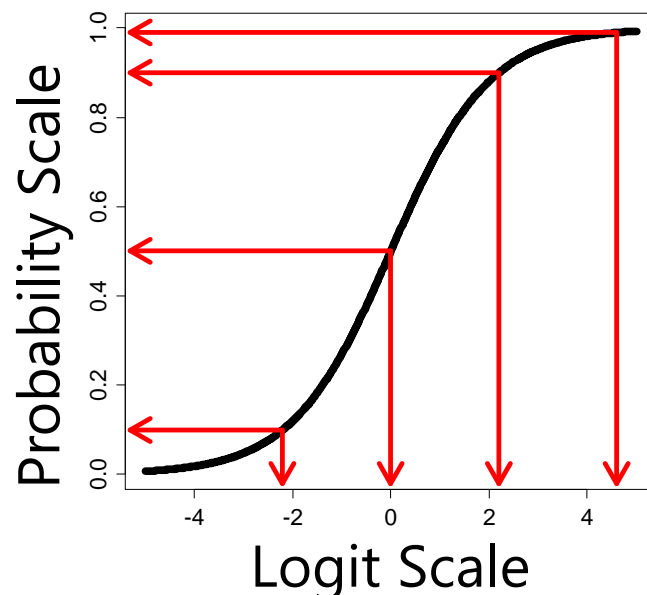
- Solution #1: Rather than predicting $P(y_i = 1)$ directly, we must transform it into an unbounded variable with a **link function**:

➤ Transform **probability** into an **odds ratio**: $\frac{p}{1-p} = \frac{\text{prob}(y=1)}{\text{prob}(y=0)}$

- If $P(y_i = 1) = .7$ then Odds(1) = 2.33; Odds(0) = .429
- But odds scale is skewed, asymmetric, and ranges from 0 to $+\infty \rightarrow$ Not helpful

➤ Take **natural log of odds ratio** \rightarrow called “**logit**” link: **Log** $\left[\frac{p}{1-p} \right]$

- If $P(y_i = 1) = .7$, then Logit(1) = .846; Logit(0) = $-.846$
- Logit scale is now symmetric about 0, range is $\pm\infty \rightarrow$ DING

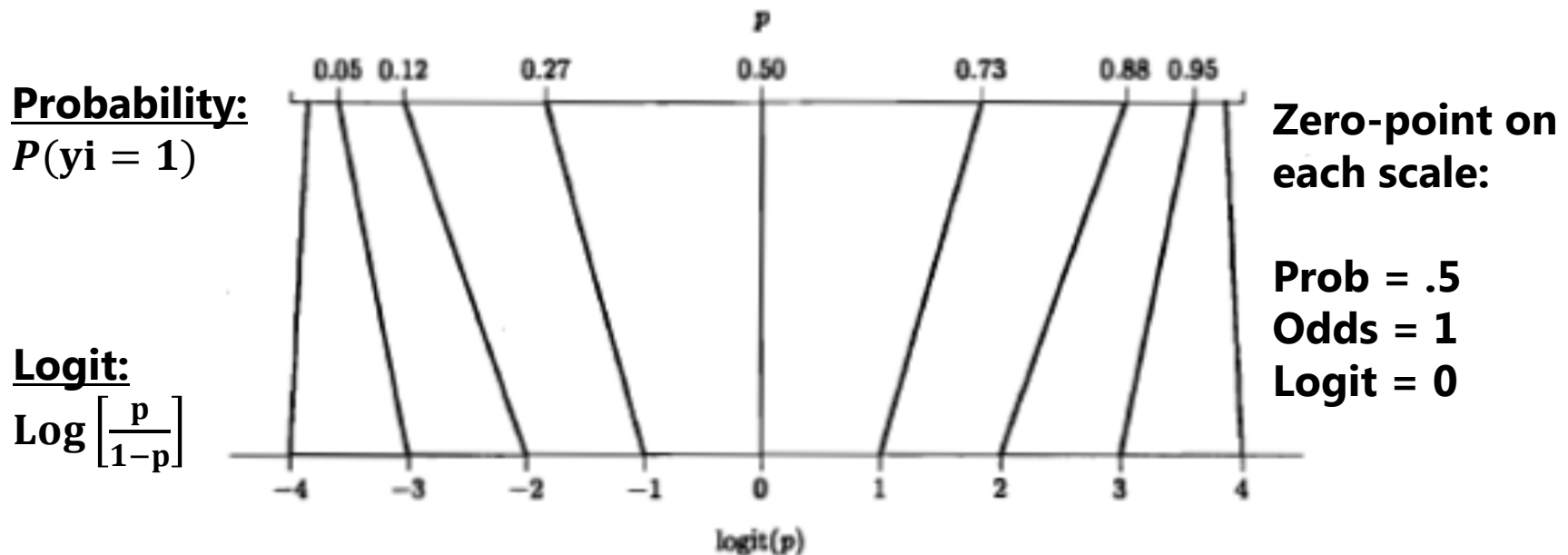


Probability	Logit
0.99	4.6
0.90	2.2
0.50	0.0
0.10	-2.2

Can you guess what $P(.01)$ would be on the logit scale?

Solution #1: Probability into Logits

- **A Logit link is a nonlinear transformation of probability:**
 - Equal intervals in logits are NOT equal intervals of probability
 - The logit goes from $\pm\infty$ and is symmetric about prob = .5 (logit = 0)
 - Now we can use a linear model → The model will be **linear with respect to the predicted logit**, which translates into a nonlinear prediction with respect to probability → **the predicted outcome shuts off at 0 or 1 as needed**



Normal GLM for Binary Outcomes?

- General linear model: $P(y_i = 1) = \beta_0 + \beta_1 X_i + \beta_2 Z_i + e_i$
- If y_i is binary, then e_i can only be 2 things: $e_i = y_i - \hat{y}_i$
 - If $y_i = 0$ then $e_i = (0 - \text{predicted probability})$
 - If $y_i = 1$ then $e_i = (1 - \text{predicted probability})$
- Problem #2a: So the residuals can't be normally distributed
- Problem #2b: The residual variance can't be constant over X as in GLM because the **mean and variance are dependent**
 - Variance of binary variable: $\text{Var}(y_i) = p * (1 - p)$

Mean and Variance of a Binary Variable

Mean (p)	.0	.1	.2	.3	.4	.5	.6	.7	.8	.9	1.0
Variance	.0	.09	.16	.21	.24	.25	.24	.21	.16	.09	.0

Solution #2: Bernoulli Distribution

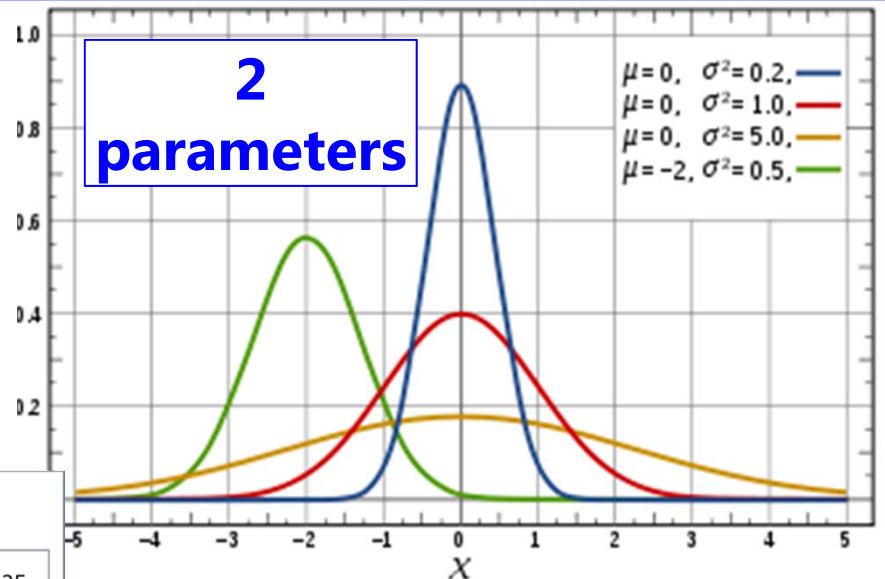
- Rather than using a normal distribution for our residuals, we will use a **Bernoulli distribution** → a special case of a binomial distribution for only one binary outcome

Univariate Normal PDF:

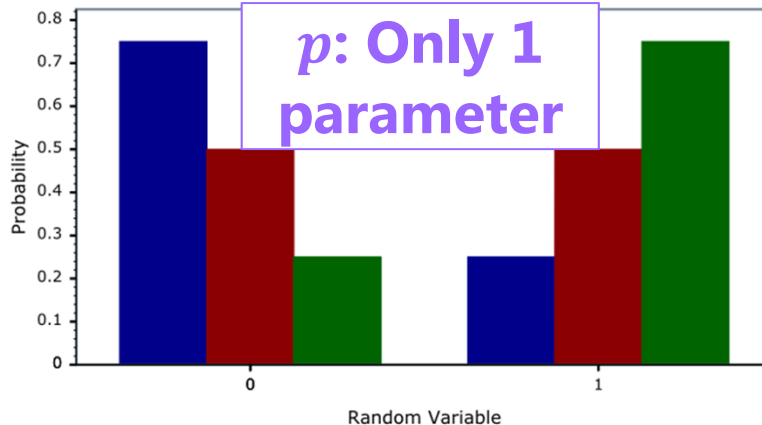
$$f(y_i) = \frac{1}{\sqrt{2\pi\sigma_e^2}} * \exp\left[-\frac{1}{2} * \frac{(y_i - \hat{y}_i)^2}{\sigma_e^2}\right]$$

Likelihood (y_i)

**2
parameters**



Bernoulli Distribution PDF



**p : Only 1
parameter**

Bernoulli PDF:

$$f(y_i) = (p_i)^{y_i} (1 - p_i)^{1-y_i}$$

**= $p(1)$ if 1,
 $p(0)$ if 0**

Predicted Binary Outcomes

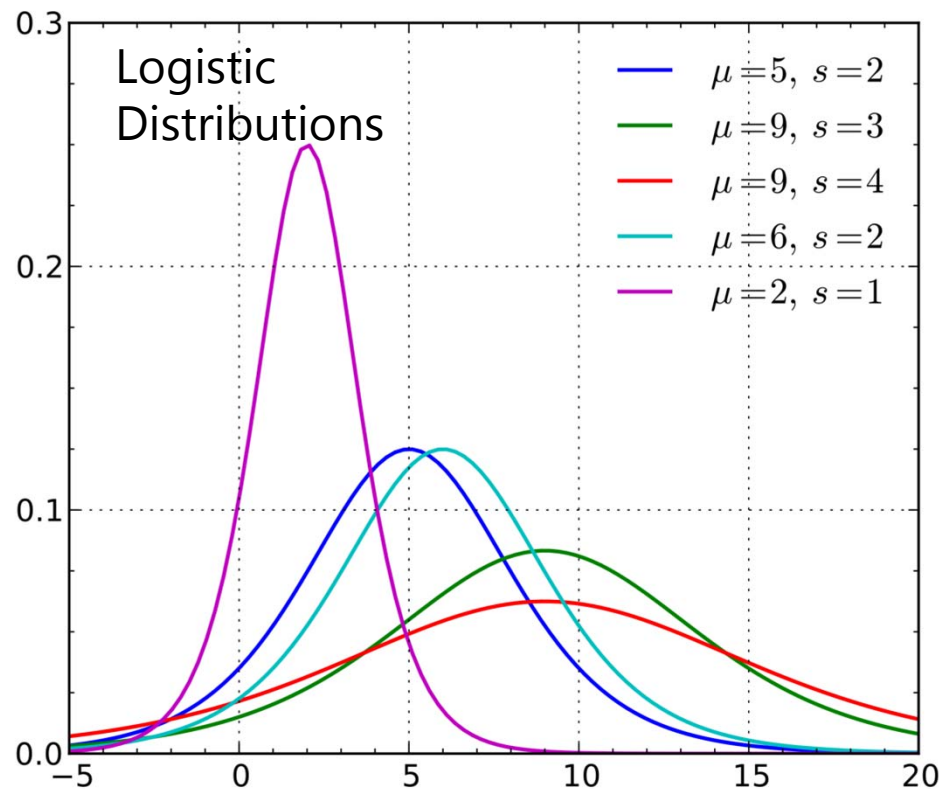
- **Logit:** $\text{Log} \left[\frac{p}{1-p} \right] = \beta_0 + \beta_1 X_i + \beta_2 Z_i$ ← g(·) link
 - Predictor effects are linear and additive like in GLM, but β = change in **logit(y)** per one-unit change in predictor
- **Odds:** $\left[\frac{p}{1-p} \right] = \exp(\beta_0) * (\beta_1 X_i) * (\beta_2 Z_i)$
or $\left[\frac{p}{1-p} \right] = \exp(\beta_0 + \beta_1 X_i + \beta_2 Z_i)$
- **Probability:** $P(y_i = 1) = \frac{\exp(\beta_0 + \beta_1 X_i + \beta_2 Z_i)}{1 + \exp(\beta_0 + \beta_1 X_i + \beta_2 Z_i)}$ ← $g^{-1}(\cdot)$
inverse
link
or $P(y_i = 1) = \frac{1}{1 + \exp[-1(\beta_0 + \beta_1 X_i + \beta_2 Z_i)]}$

“Logistic Regression” for Binary Data

- This model is sometimes expressed by calling the $\text{logit}(y_i)$ a underlying continuous (“latent”) response of y_i^* instead:

$$y_i^* = \text{threshold} + \text{your model} + e_i \quad \text{threshold} = \beta_0 * -1$$

- In which $y_i = 1$ if $(y_i^* > \text{threshold})$, or $y_i = 0$ if $(y_i^* \leq \text{threshold})$



So **if predicting y_i^*** , then

$$e_i \sim \text{Logistic}(0, \sigma_e^2 = 3.29)$$

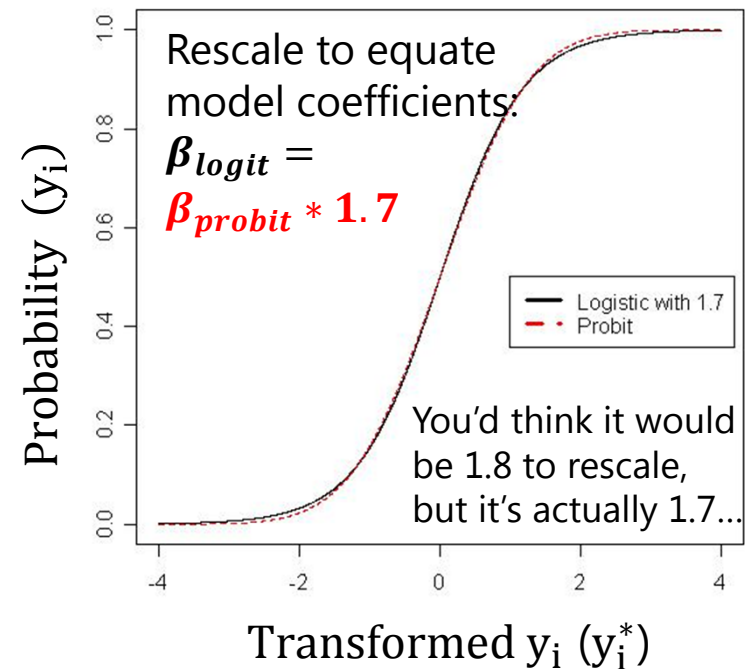
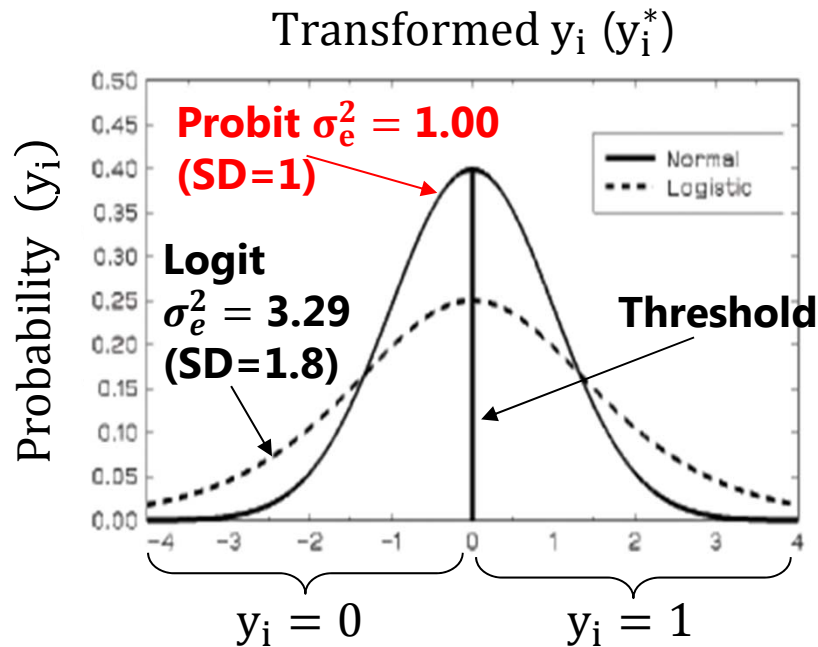
Logistic Distribution:

Mean = μ , Variance = $\frac{\pi^2}{3} s^2$,
where s = scale factor that
allows for “over-dispersion”
(must be fixed to 1 in logistic
regression for identification)

Other Models for Binary Data

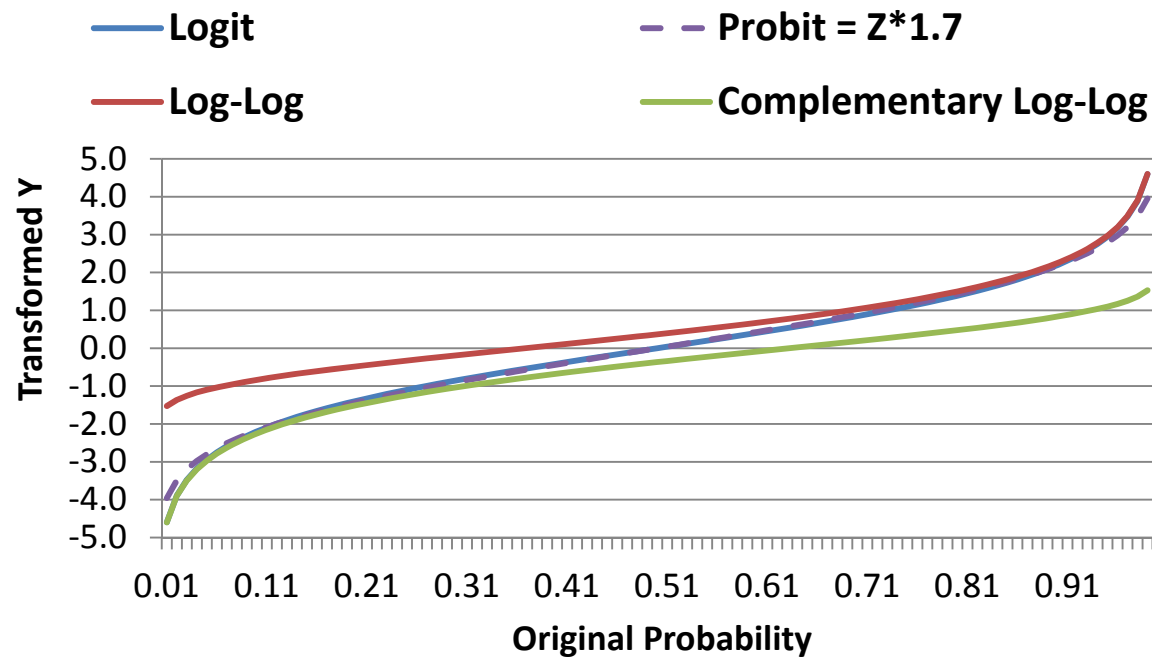
- The idea that a “latent” continuous variable underlies an observed binary response also appears in a **Probit Regression** model:
 - A **probit** link, such that now your model predicts a different transformed Y_p :
$$\text{Probit}(y_i = 1) = \Phi^{-1}P(y_i = 1) = \text{your model} \quad \leftarrow \boxed{g(\cdot)}$$
 - Where Φ = standard normal cumulative distribution function, so the transformed y_i is the **z-score** that corresponds to the value of standard normal curve below which observed probability is found (requires integration to transform back)
 - Same binomial (Bernoulli) distribution for the binary e_i residuals, in which residual variance cannot be separately estimated (so no e_i in the model)
 - Probit also predicts “latent” response: $y_i^* = \text{threshold} + \text{your model} + e_i$
 - But Probit says $e_i \sim \text{Normal}(0, \sigma_e^2 = 1.00)$, whereas Logit $\sigma_e^2 = \frac{\pi^2}{3} = 3.29$
 - So given this difference in variance, probit estimates are on a different scale than logit estimates, and so their estimates won’t match... however...

Probit vs. Logit: Should you care? Pry not.



- Other fun facts about probit:
 - Probit = "ogive" in the Item Response Theory (IRT) world
 - Probit has no odds ratios (because it's not based on odds)
- Both logit and probit assume **symmetry** of the probability curve, but there are other *asymmetric* options as well...

Other Link Functions for Binary Outcomes



Logit = Probit*1.7
which both assume
symmetry of prediction

**Log-Log is for outcomes in
which 1 is more frequent**

**Complementary
Log-Log is for outcomes in
which 0 is more frequent**

$\mu = \text{model}$	Logit	Probit	Log-Log	Complement. Log-Log
$g(\cdot)$ for new y_i :	$\text{Log}\left(\frac{p}{1-p}\right) = \mu$	$\Phi^{-1}(p) = \mu$	$-\text{Log}[-\text{Log}(p)] = \mu$	$\text{Log}[-\text{Log}(1-p)] = \mu$
$g^{-1}(\cdot)$ to get back to probability:	$p = \frac{\exp(\mu)}{1 + \exp(\mu)}$	$p = \Phi(\mu)$	$p = \exp[-\exp(-\mu)]$ $e_i \sim \text{extreme value} \left(-\gamma?, \sigma_e^2 = \frac{\pi^2}{6} \right)$	$p = 1 - \exp[-\exp(\mu)]$
In SAS LINK=	LOGIT	PROBIT	LOGLOG	CLOGLOG

Generalized Multilevel Models for Non-Normal Longitudinal Data

- Topics:
 - 3 parts of a generalized (multilevel) model
 - Models for binary outcomes
 - **Models for categorical outcomes**
 - Complications for generalized multilevel models
 - A brief tour of other generalized models:
 - Models for count outcomes
 - Models for not-normal but continuous outcomes

Too Logit to Quit... <http://www.youtube.com/watch?v=Cdk1gwWH-Cg>

- The **logit** is the basis for many other generalized models for predicting categorical outcomes
- Next we'll see how C possible response categories can be predicted using $C - 1$ binary "submodels" that involve carving up the categories in different ways, in which each binary submodel uses a logit link to predict its outcome
- Types of categorical outcomes:
 - Definitely ordered categories: "**cumulative logit**"
 - Maybe ordered categories: "**adjacent category logit**" (not used much)
 - Definitely NOT ordered categories: "**generalized logit**"

Logit-Based Models for C Ordinal Categories

- Known as “**cumulative logit**” or “**proportional odds**” model in generalized models; known as “graded response model” in IRT
 - LINK=CLOGIT, (DIST=MULT) in SAS GLIMMIX
- Models the probability of **lower vs. higher** cumulative categories via $C - 1$ submodels (e.g., if $C = 4$ possible responses of $c = 0,1,2,3$):

0 vs. **1, 2, 3**
Submodel₁

0,1 vs. **2,3**
Submodel₂

0,1,2 vs. **3**
Submodel₃

I've named these submodels based on what they predict, but SAS will name them its own way in the output.

- In SAS, what the binary submodels predict depends on whether the model is predicting **DOWN** ($y_i = 0$, the default) or **UP** ($y_i = 1$) **cumulatively**
- Example predicting UP in an empty model (subscripts=parm,submodel)**
- Submodel 1: $\text{Logit}(y_i > 0) = \beta_{01} \rightarrow P(y_i > 0) = \exp(\beta_{01})/[1 + \exp(\beta_{01})]$
- Submodel 2: $\text{Logit}(y_i > 1) = \beta_{02} \rightarrow P(y_i > 1) = \exp(\beta_{02})/[1 + \exp(\beta_{02})]$
- Submodel 3: $\text{Logit}(y_i > 2) = \beta_{03} \rightarrow P(y_i > 2) = \exp(\beta_{03})/[1 + \exp(\beta_{03})]$

Logit-Based Models for C Ordinal Categories

- Models the probability of **lower vs. higher** cumulative categories via $C - 1$ submodels (e.g., if $C = 4$ possible responses of $c = 0,1,2,3$):

0 vs. **1,2,3**
Submodel₁
→ Prob₁

0,1 vs. **2,3**
Submodel₂
→ Prob₂

0,1,2 vs. **3**
Submodel₃
→ Prob₃

$$\text{Logit}(y_i > 2) = \beta_{03}$$

$$\rightarrow P(y_i > 2) = \frac{\exp(\beta_{03})}{1 + \exp(\beta_{03})}$$

- In SAS, what the binary submodels predict depends on whether the model is predicting **DOWN** ($y_i = 0$, the default) or **UP** ($y_i = 1$) **cumulatively**
 - Either way, the model predicts the middle category responses *indirectly*
- Example if predicting UP with an empty model:**

- Probability of 0 = $1 - \text{Prob}_1$
- Probability of 1 = $\text{Prob}_1 - \text{Prob}_2$
- Probability of 2 = $\text{Prob}_2 - \text{Prob}_3$
- Probability of 3 = $\text{Prob}_3 - 0$

The cumulative submodels that create these probabilities are each estimated using **all the data** (good, especially for categories not chosen often), but **assume order in doing so** (maybe bad, maybe ok, depending on your response format).

Logit-Based Models for C Ordinal Categories

- Ordinal models usually use a logit link transformation, but they can also use cumulative log-log or cumulative complementary log-log links
 - LINK= CUMLOGLOG or CUMCLL, respectively, in SAS PROC GLIMMIX
- Almost always assume **proportional odds**, that effects of predictors are the same across binary submodels—for example (subscripts = parm, submodel)
 - Submodel 1: $\text{Logit}(y_i > 0) = \beta_{01} + \beta_1 X_i + \beta_2 Z_i + \beta_3 X_i Z_i$
 - Submodel 2: $\text{Logit}(y_i > 1) = \beta_{02} + \beta_1 X_i + \beta_2 Z_i + \beta_3 X_i Z_i$
 - Submodel 3: $\text{Logit}(y_i > 2) = \beta_{03} + \beta_1 X_i + \beta_2 Z_i + \beta_3 X_i Z_i$
- Proportional odds essentially means no interaction between submodel and predictor effects, which greatly reduces the number of estimated parameters
 - Assumption for single-level data can be tested painlessly using PROC LOGISTIC, which provides a global SCORE test of equivalence of all slopes between submodels
 - If the proportional odds assumption fails and $C > 3$, you'll need to write your own model non-proportional odds ordinal model in PROC NLMIXED

Logit-Based Models for C Categories

- Uses multinomial distribution, whose PDF for $C = 4$ categories of $c = 0, 1, 2, 3$, an observed $y_i = c$, and indicators I if $c = y_i$

$$f(y_i = c) = p_{i0}^{I[y_i=0]} p_{i1}^{I[y_i=1]} p_{i2}^{I[y_i=2]} p_{i3}^{I[y_i=3]}$$

Only p_{ic} for the response $y_i = c$ gets used

- Maximum likelihood is then used to find the most likely parameters in the model to predict the probability of each response through the (usually logit) link function; probabilities sum to 1: $\sum_{c=1}^C p_{ic} = 1$
- Other models for categorical data that use the multinomial:
 - Adjacent category logit (partial credit): Models the probability of **each next highest** category via $C - 1$ submodels (e.g., if $C = 4$):

0 vs. 1 1 vs. 2 2 vs. 3
 - Baseline category logit (nominal): Models the probability of **reference vs. other** category via $C - 1$ submodels (e.g., if $C = 4$ and $0 = \text{ref}$):

0 vs. 1 0 vs. 2 0 vs. 3

Generalized Multilevel Models for Non-Normal Longitudinal Data

- Topics:
 - 3 parts of a generalized (multilevel) model
 - Models for binary outcomes
 - Models for categorical outcomes
 - **Complications for generalized multilevel models**
 - A brief tour of other generalized models:
 - Models for count outcomes
 - Models for not-normal but continuous outcomes

From Single-Level to Multilevel...

- Multilevel generalized models have the same 3 parts as single-level generalized models:
 - Link function to transform bounded DV into unbounded outcome
 - Linear model that directly predicts link-transformed DV instead
 - Alternative distribution of level-1 residuals used (e.g., Bernoulli)
- But in adding random effects (i.e., additional piles of variance) to address dependency in longitudinal data:
 - Piles of variance are ADDED TO, not EXTRACTED FROM, the original residual variance pile when it is fixed to a known value (e.g., 3.29), which causes the model coefficients to change scale across models
 - ML estimation is way more difficult because normal random effects + not-normal residuals does not have a known distribution like MVN
 - No such thing as REML for generalized multilevel models

Empty Multilevel Model for Binary Outcomes

- **Level 1:** $\text{Logit}(y_{ti}) = \beta_{0i}$
- **Level 2:** $\beta_{0i} = \gamma_{00} + u_{0i}$
- **Composite:** $\text{Logit}(y_{ti}) = \gamma_{00} + u_{0i}$
- σ_e^2 residual variance is not estimated $\rightarrow \pi^2/3 = 3.29$
 - (Known) residual is in model for actual Y, not prob(Y) or logit(Y)
- Logistic ICC = $\frac{BP}{BP+WP} = \frac{\tau_{U0}^2}{\tau_{U0}^2 + \sigma_e^2} = \frac{\tau_{U0}^2}{\tau_{U0}^2 + 3.29}$
- Can do $-2\Delta LL$ test to see if $\tau_{U0}^2 > 0$, although the ICC is somewhat problematic to interpret due to non-constant residual variance

Note what's
NOT in level 1...

Random Linear Time Model for Binary Outcomes

- **Level 1:** $\text{Logit}(y_{ti}) = \beta_{0i} + \beta_{1i}(\text{time}_{ti})$
- **Level 2:** $\beta_{0i} = \gamma_{00} + u_{0i}$
 $\beta_{1i} = \gamma_{10} + u_{1i}$
- **Combined:** $\text{Logit}(y_{ti}) = (\gamma_{00} + u_{0i}) + (\gamma_{10} + u_{1i})(\text{time}_{ti})$
- σ_e^2 residual variance is still not estimated $\rightarrow \pi^2/3 = 3.29$
- Can test new fixed or random effects with $-2\Delta\text{LL}$ tests
(or Wald test p -values for fixed effects as usual)

Random Linear Time Model for **Ordinal** Outcomes ($C = 3$)

- **L1:** $\text{Logit}(y_{ti1}) = \beta_{0i1} + \beta_{1i1}(\text{time}_{ti})$
 $\text{Logit}(y_{ti2}) = \beta_{0i2} + \beta_{1i2}(\text{time}_{ti})$
- **L2:** $\beta_{0i1} = Y_{001} + U_{0i1}$ $\beta_{1i1} = Y_{101} + U_{1i1}$
 $\beta_{0i2} = Y_{002} + U_{0i2}$ $\beta_{1i2} = Y_{102} + U_{1i2}$
- Assumes proportional odds \rightarrow
 $Y_{001} \neq Y_{002}$ and $Y_{101} = Y_{102}$ and $U_{0i1} = U_{0i2}$ and $U_{1i1} = U_{1i2}$
 - Testable via nominal model (all unequal) or using NLMIXED to write a custom model in which some can be constrained
 - σ_e^2 residual variance is still not estimated $\rightarrow \pi^2/3 = 3.29$

New Interpretation of Fixed Effects

- In general linear mixed models, the fixed effects are interpreted as the “average” effect for the sample
 - γ_{00} is “sample average” intercept
 - u_{0i} is “individual deviation from sample average”
- What “average” means in *generalized* linear mixed models is different, because the natural log is a nonlinear function:
 - So the mean of the logs \neq log of the means
 - Therefore, the fixed effects are not the “sample average” effect, they are the effect for ***specifically for $U_i = 0$***
 - Fixed effects are *conditional* on the random effects
 - This gets called a “unit-specific” or “subject-specific” model
 - This distinction does not exist for normally distributed outcomes

Comparing Results across Models

- NEW RULE: Coefficients cannot be compared across models, because they are not on the same scale! (see Bauer, 2009)
- e.g., if residual variance = 3.29 in binary models:
 - When adding a random intercept variance to an empty model, the **total variation in the outcome has increased** → the fixed effects will increase in size because they are *unstandardized* slopes

$$\gamma_{\text{mixed}} \approx \sqrt{\frac{\tau_{U_0}^2 + 3.29}{3.29}} (\beta_{\text{fixed}})$$

- **Level-1 predictors cannot decrease the residual variance** like usual, so all other models estimates have to go up to compensate
 - If X_{ti} is uncorrelated with other X 's and is a pure level-1 variable ($\text{ICC} \approx 0$), then fixed and $\text{SD}(U_{0i})$ will increase by same factor
- **Random effects variances can decrease**, though, so level-2 effects should be on the same scale across models if level-1 is the same

A Little Bit about Estimation

- Goal: End up with maximum likelihood estimates for all model parameters (because they are consistent, efficient)
 - When we have a **V** matrix based on multivariate **normally** distributed \mathbf{e}_{ti} residuals at level-1 and multivariate normally distributed \mathbf{U}_i terms at level 2, ML is easy
 - When we have a **V** matrix based on multivariate **Bernoulli** distributed \mathbf{e}_{ti} residuals at level-1 and multivariate normally distributed \mathbf{U}_i terms at level 2, ML is much harder
 - Same with any other kind model for “not normal” level 1 residual
 - **ML does not assume normality unless you fit a “normal” model!**
- 3 main families of estimation approaches:
 - Quasi-Likelihood methods (“marginal/penalized quasi ML”)
 - Numerical Integration (“adaptive Gaussian quadrature”)
 - Also Bayesian methods (MCMC, newly available in SAS or Mplus)

2 Main Types of Estimation

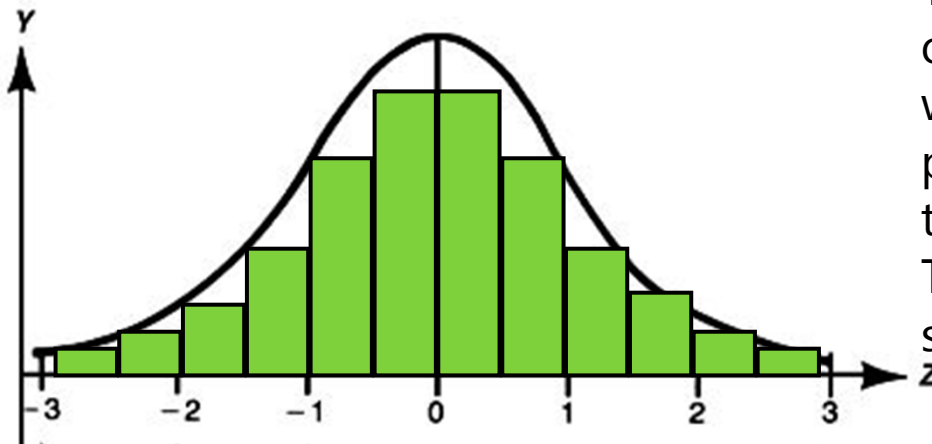
- **Quasi-Likelihood methods** → older methods
 - “Marginal QL” → approximation around fixed part of model
 - “Penalized QL” → approximation around fixed + random parts
 - These both underestimate variances (MQL more so than PQL)
 - 2nd-order PQL is supposed to be better than 1st-order MQL
 - QL methods DO NOT PERMIT MODEL $-2\Delta LL$ TESTS
 - HLM program adds Laplace approximation to QL, which then does permit $-2\Delta LL$ tests (also in SAS GLIMMIX and STATA xtmelogit)
- **ML via Numerical Integration** → gold standard
 - Much better estimates and $-2\Delta LL$ tests, but can take for-freaking-ever (can use PQL methods to get good start values)
 - Will blow up with many random effects (which make the model exponentially more complex, especially in these models)
 - Relies on assumptions of local independence, like usual → all level-1 dependency has been modeled; level-2 units are independent

ML via Numerical Integration

- **Step 1:** Select **starting values** for all fixed effects
- **Step 2:** Compute the **likelihood** of each observation given by the *current* parameter values using chosen distribution of residuals
 - Model gives link-predicted outcome given parameter estimates, but the U 's themselves are not parameters—their variance is instead
 - But so long as we can assume the U 's are MVN, we can still proceed
 - Computing the likelihood for each set of possible parameters requires *removing* the individual U values from the model equation—by **integrating** across possible U values for each Level-2 unit
 - Integration is accomplished by “Gaussian Quadrature” → summing up rectangles that approximate the integral (area under the curve) for each Level-2 unit
- **Step 3:** Decide if you have the right answers, which occurs when the log-likelihood changes very little across iterations (i.e., it converges)
- **Step 4:** If you aren't converged, choose new parameters values
 - Newton-Raphson or Fisher Scoring (calculus), EM algorithm (U 's = missing data)

ML via Numerical Integration

- More on Step 2: Divide the U distribution into rectangles
 - → “Gaussian Quadrature” (# rectangles = # “quadrature points”)
 - Can either divide the whole distribution into rectangles, or take the most likely section for each level-2 unit and rectangle that
 - This is “adaptive quadrature” and is computationally more demanding, but gives more accurate results with fewer rectangles



The likelihood of each level-2 unit's outcomes at each **U** rectangle is then weighted by that rectangle's probability of being observed (from the multivariate normal distribution). The weighted likelihoods are then summed across all rectangles...

→ ta da! “**numerical integration**”

Example of Numeric Integration: Binary DV, Fixed Linear Time, Random Intercept Model

1. Start with values for fixed effects: intercept: $\gamma_{00} = 0.5$, time: $\gamma_{10} = 2.0$,
2. Compute likelihood for real data based on fixed effects and plausible U_{0i} (-2,0,2) using model: $\text{Logit}(y_{ti}=1) = \gamma_{00} + \gamma_{10}(\text{time}_{ti}) + U_{0i}$
 - Here for one person at two occasions with $y_{ti}=1$ at both occasions

			IF $y_{ti}=1$	IF $y_{ti}=0$	Likelihood	Theta	Theta	Product
	$U_{0i} = -2$	$\text{Logit}(y_{ti})$	Prob	1-Prob	if both $y=1$	prob	width	per Theta
Time 0	$0.5 + 1.5(0) - 2$	-1.5	0.18	0.82	0.091213	0.05	2	0.00912
Time 1	$0.5 + 1.5(1) - 2$	0.0	0.50	0.50				
	$U_{0i} = 0$	$\text{Logit}(y_{ti})$	Prob	1-Prob				
Time 0	$0.5 + 1.5(0) + 0$	0.5	0.62	0.38	0.54826	0.40	2	0.43861
Time 1	$0.5 + 1.5(1) + 0$	2.0	0.88	0.12				
	$U_{0i} = 2$	$\text{Logit}(y_{ti})$	Prob	1-Prob				
Time 0	$0.5 + 1.5(0) + 2$	2.5	0.92	0.08	0.90752	0.05	2	0.09075
Time 1	$0.5 + 1.5(1) + 2$	4.0	0.98	0.02				
Overall Likelihood (Sum of Products over All Thetas):								0.53848

(do this for each occasion, then multiply this whole thing over all people)

(repeat with new values of fixed effects until find highest overall likelihood)

Summary: Generalized Multilevel Models

- Analyze link-transformed DV (e.g., via logit, log, log-log...)
 - **Linear** relationship between X's and **transformed** Y
 - **Nonlinear** relationship between X's and **original** Y
 - Original e_{ti} residuals are assumed to follow some non-normal distribution
- In models for binary or categorical data, Level-1 residual variance is set
 - So it can't go down after adding level-1 predictors, which means that the scale of everything else has to go UP to compensate
 - Scale of model will also be different after adding random effects for the same reason—the total variation in the model is now bigger
 - Fixed effects may not be comparable across models as a result
- Estimation is trickier and takes longer
 - Numerical integration is best but may blow up in complex models
 - Start values are often essential (can get those with MSPL estimator)

Generalized Multilevel Models for Non-Normal Longitudinal Data

- Topics:
 - 3 parts of a generalized (multilevel) model
 - Models for binary outcomes
 - Models for categorical outcomes
 - Complications for generalized multilevel models
 - **A brief tour of other generalized models:**
 - **Models for count outcomes**
 - **Models for not-normal but continuous outcomes**

A Taxonomy of Not-Normal Outcomes

- **“Discrete” outcomes**—all responses are **whole** numbers
 - **Categorical variables** in which **values are labels**, not amounts
 - Binomial (2 options) or multinomial (3+ options) distributions
 - Question: Are the values ordered → which link?
 - **Count of things that happened**, so values < 0 cannot exist
 - Sample space goes from 0 to positive infinity
 - Poisson or Negative Binomial distributions (usually)
 - Log link (usually) so predicted outcomes can't go below 0
 - Question: Are there *extra* 0 values? What to do about them?
- **“Continuous” outcomes**—responses can be **any** number
 - Question: What does the residual distribution look like?
 - Normal-ish? Skewed? Cut off? Mixture of different distributions?

Models for Count Outcomes

- Counts: non-negative integer unbounded responses
 - e.g., how many cigarettes did you smoke this week?
 - Traditionally uses natural log link so that predicted outcomes stay ≥ 0
- $g(\bullet)$ $\text{Log}(E(y_i)) = \text{Log}(\mu_i) = \text{model} \rightarrow$ predicts mean of y_i
- $g^{-1}(\bullet)$ $E(y_i) = \exp(\text{model}) \rightarrow$ to un-log it, use $\exp(\text{model})$
 - e.g., if $\text{Log}(\mu_i) = \text{model}$ provides predicted $\text{Log}(\mu_i) = 1.098$, that translates to an actual predicted count of $\exp(1.098) = 3$
 - e.g., if $\text{Log}(\mu_i) = \text{model}$ provides predicted $\text{Log}(\mu_i) = -5$, that translates to an actual predicted count of $\exp(-5) = 0.006738$
- So that's how linear model predicts μ_i , the expected count for y_i , but what about residual variance?

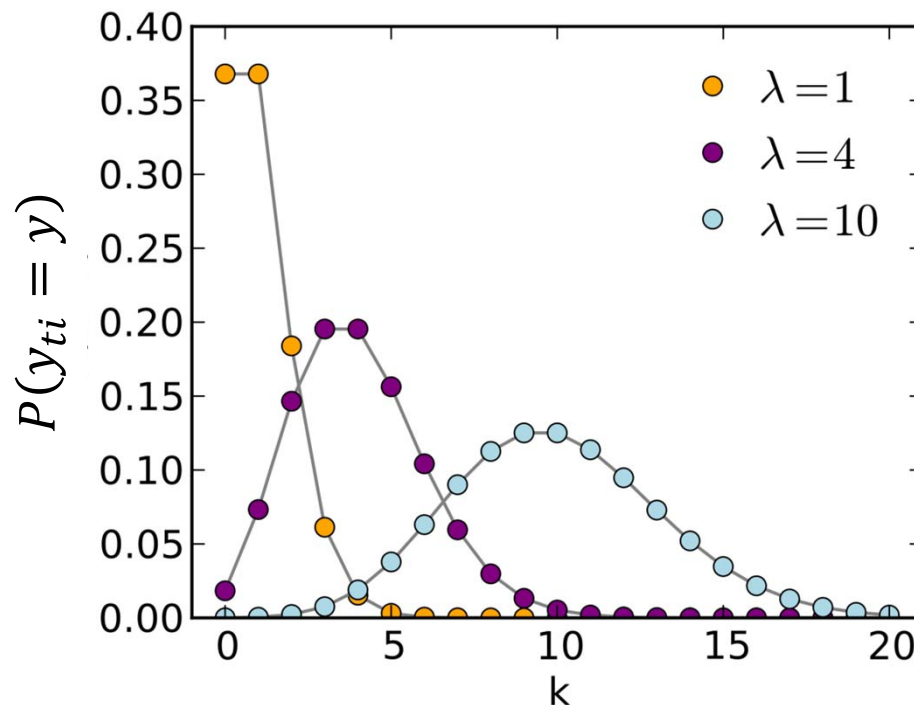
Poisson Distribution for Residuals

- Poisson distribution has one parameter, λ , which is both its mean and its variance (so $\lambda = \text{mean} = \text{variance}$ in Poisson)

- $f(y_i|\lambda) = \text{Prob}(y_i = y) = \frac{\lambda^y * \exp(-\lambda)}{y!}$

$y!$ is factorial of y

- PDF: $\text{Prob}(y_i = y|\beta_0, \beta_1, \beta_2) = \frac{\mu_i^y * \exp(-\mu_i)}{y!}$



The dots indicate that only integer values are observed.

Distributions with a small expected value (mean or λ) are predicted to have a lot of 0's.

Once $\lambda > 6$ or so, the shape of the distribution is close to that of a normal distribution.

3 potential problems for Poisson...

- The standard Poisson distribution is rarely sufficient, though
- **Problem #1: When mean \neq variance**
 - If variance < mean, this leads to “under-dispersion” (not that likely)
 - If variance > mean, this leads to “over-dispersion” (happens frequently)
- **Problem #2: When there are *no* 0 values**
 - Some 0 values are expected from count models, but in some contexts $y_i > 0$ always (but subtracting 1 won't fix it; need to adjust the model)
- **Problem #3: When there are *too many* 0 values**
 - Some 0 values are expected from the Poisson and Negative Binomial models already, but many times there are even more 0 values observed than that
 - To fix it, there are two main options, depending on what you do to the 0's
- Each of these problems requires a model adjustment to fix it...

Problem #1: Variance > mean = over-dispersion

- To fix it, we must add another parameter that allows the variance to exceed the mean... becomes a Negative Binomial distribution
 - Says residuals are a mixture of Poisson and gamma distributions

- Model: $\text{Log}(y_i) = \text{Log}(\mu_i) = \beta_0 + \beta_1 X_i + \beta_2 Z_i + e_i^G$

- Poisson PDF was: $\text{Prob}(y_i = y | \beta_0, \beta_1, \beta_2) = \frac{\mu_i^y \cdot \exp(-\mu_i)}{y!}$

- Negative Binomial PDF with a new k dispersion parameter is now:

- $\text{Prob}(y_i = y | \beta_0, \beta_1, \beta_2) = \frac{\Gamma(y + \frac{1}{k})}{\Gamma(y+1) \cdot \Gamma(\frac{1}{k})} * \frac{(k\mu_i)^y}{(1+k\mu_i)^{y+\frac{1}{k}}}$

**DIST =
NEGBIN** in SAS

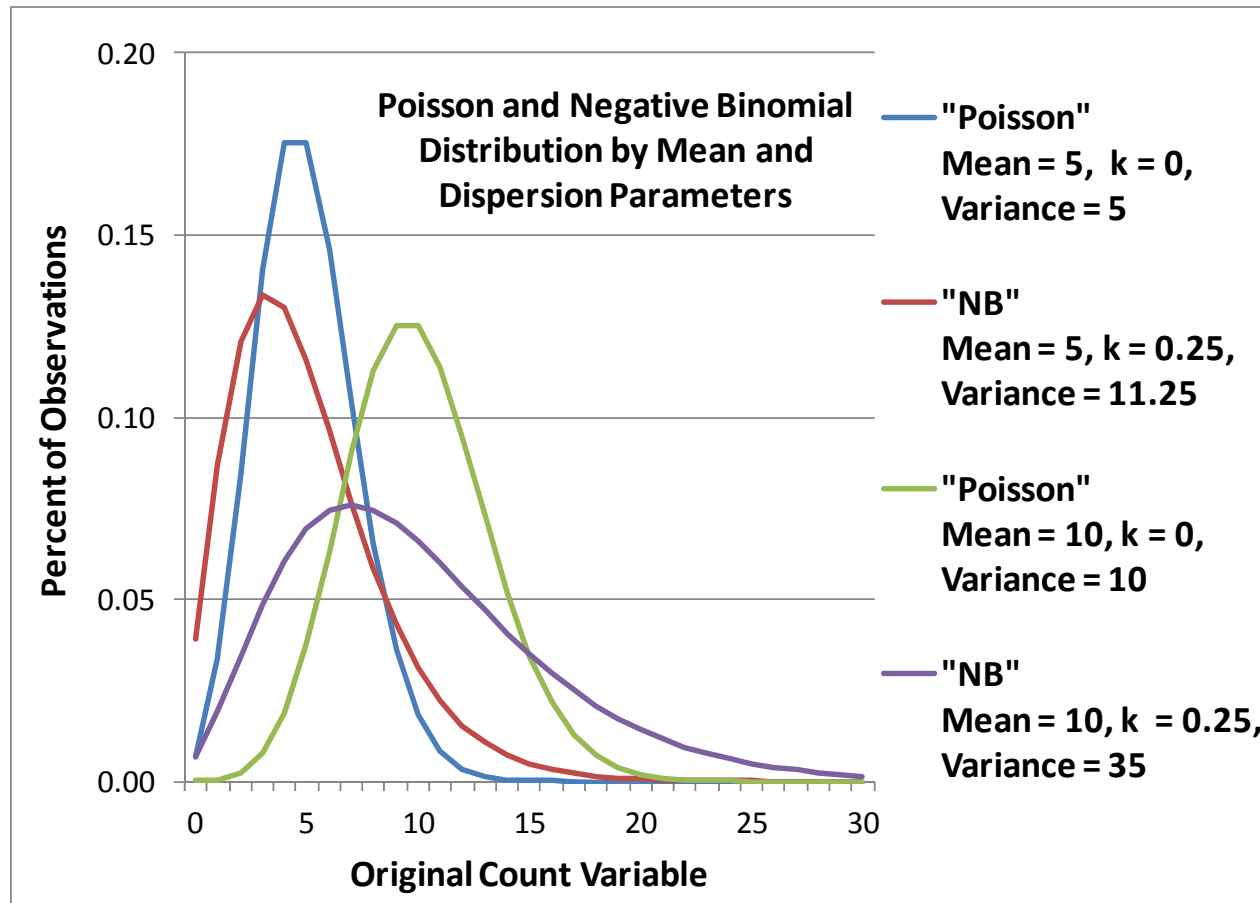
- k is dispersion, such that $\text{Var}(y_i) = \mu_i + k\mu_i^2$

So is Poisson if $k = 0$

- Non-Poisson related $e_i^G \sim \text{Gamma}(\text{mean} = 1, \text{variance} = k)$

- Since $\text{Log}(1) = 0$, the extra 0's won't add to the predicted log count, and if there is no extra dispersion, then variance of $e_i^G \sim 0$

Negative Binomial (NB) = “Stretchy” Poisson...



Mean = λ

Dispersion = k

$$\text{Var}(y_i) = \lambda + k\lambda^2$$

A Negative Binomial model can be useful for count outcome residuals that have some extra skewness, but otherwise follow a Poisson distribution.

- Because its k dispersion parameter is fixed to 0, the Poisson model is nested within the Negative Binomial model—to test improvement in fit:
- Is $-2(LL_{\text{Poisson}} - LL_{\text{NegBin}}) > 3.84$ for $df = 1$? Then $p < .05$, keep NB

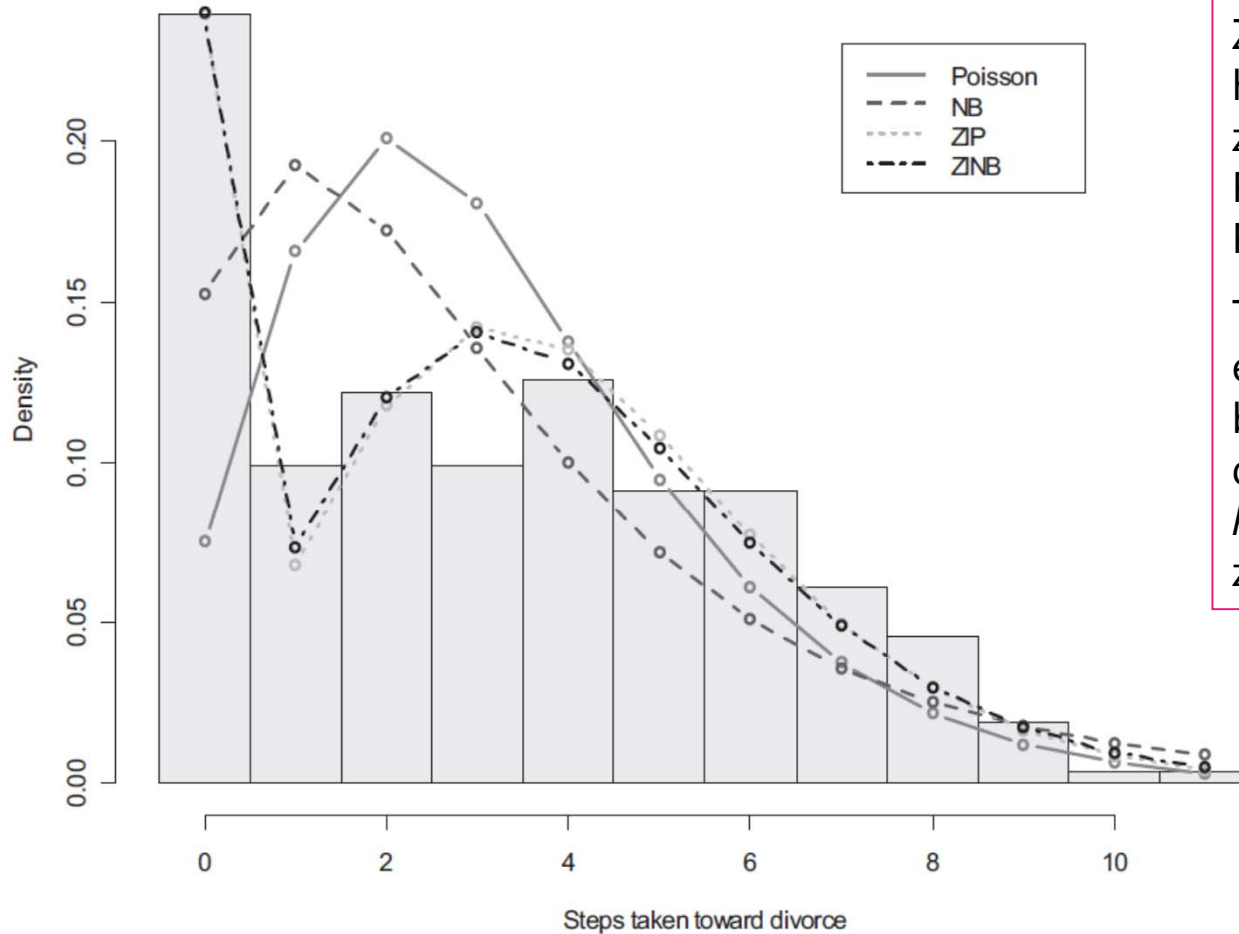
Problem #2: There are no 0 values

- **"Zero-Altered"** or **"Zero-Truncated"** Poisson or Negative Binomial: ZAP/ZANB or ZTP/ZTNB
 - Is usual count distribution, just not allowing any 0 values
 - Poisson version is readily available within SAS PROC FMM using DIST=TRUNCPOISSON (next version should have TRUNCNEGBIN, too)
- Poisson PDF was: $\text{Prob}(y_i = y | \mu_i) = \frac{\mu_i^y \exp(-\mu_i)}{y!}$
- Zero-Truncated Poisson PDF is:
 - $\text{Prob}(y_i = y | \mu_i, y_i > 0) = \frac{\mu_i^y \exp(-\mu_i)}{y! [1 - \exp(-\mu_i)]}$
 - $\text{Prob}(y_i = 0) = \exp(-\mu_i)$, so $\text{Prob}(y_i > 0) = 1 - \exp(-\mu_i)$
 - Divides by probability of non-0 outcomes so probability still sums to 1

Problem #3: Too many 0 values, Option #1

- **"Zero-Inflated"** Poisson (DIST=ZIP) or Negative Binomial (DIST=ZINB); available within SAS PROC GENMOD (and Mplus)
 - Distinguishes **two kinds of 0 values: expected** and **inflated** ("structural") through a mixture of distributions (Bernoulli + Poisson/NB)
 - Creates two submodels to predict "if *extra* 0" and "if not, how much"?
 - Does not readily map onto most hypotheses (in my opinion)
 - But a ZIP example would look like this... (ZINB would add k dispersion, too)
- Submodel 1: $\text{Logit}(y_i = \text{extra } 0) = \beta_{01} + \beta_{11}X_i + \beta_{21}Z_i$
 - Predict being an extra 0 using Link = Logit, Distribution = Bernoulli
 - Don't have to specify predictors for this part, can simply allow an intercept (but need ZEROMODEL option to include predictors in SAS GENMOD)
- Submodel 2: $\text{Log}(E(y_i)) = \beta_{02} + \beta_{12}X_i + \beta_{22}Z_i$
 - Predict rest of counts (including 0's) using Link = Log, Distribution = Poisson

Example of Zero-Inflated Outcomes



Zero-inflated distributions have extra "structural zeros" not expected from Poisson or NB ("stretched Poisson") distributions.

This can be tricky to estimate and interpret because the model distinguishes between *kinds of zeros* rather than zero or not...

Image borrowed from Atkins & Gallop, 2007

Figure 1. Histogram of Marital Status Inventory with predicted probabilities from regressions. NB = negative binomial; ZIP = zero-inflated Poisson; ZINB = zero-inflated negative binomial.

Problem #3: Too many 0 values, Option #1

- The Zero-Inflated models get put back together as follows:

- ω_i is the predicted probability of being an extra 0, from:

$$\omega_i = \frac{\exp[\text{Logit}(y_i = \text{extra } 0)]}{1 + \exp[\text{Logit}(y_i = \text{extra } 0)]}$$

- μ_i is the predicted count for the rest of the distribution, from:

$$\mu_i = \exp[\text{Log}(y_i)]$$

- ZIP: Mean (original y_i) = $(1 - \omega_i)\mu_i$

- ZIP: Variance(original y_i) = $\mu_i + \frac{\omega_i}{(1-\omega_i)} \mu_i^2$

- ZINB: Mean (original y_i) = $(1 - \omega_i)\mu_i$

- ZINB: Variance(original y_i) = $\mu_i + \left[\frac{\omega_i}{(1-\omega_i)} + \frac{k}{1-\omega_i} \right] \mu_i^2$

Problem #3: Too many 0 values, Option #2

- **"Hurdle"** models for Poisson or Negative Binomial
 - PH or NBH: Explicitly **separates 0 from non-0 values** through a mixture of distributions (Bernoulli + Zero-Altered Poisson/NB)
 - Creates two submodels to predict "if any 0" and "if not 0, how much?"
 - Easier to think about in terms of prediction (in my opinion)
- Submodel 1: $\text{Logit}(y_i = 0) = \beta_{01} + \beta_{11}X_i + \beta_{21}Z_i$
 - Predict being a 0 using Link = Logit, Distribution = Bernoulli
 - Don't have to specify predictors for this part, can simply allow it to exist
- Submodel 2: $\text{Log}(E(y_i) > 0) = \beta_{02} + \beta_{12}X_i + \beta_{22}Z_i$
 - Predict rest of positive counts using Link = Log, Distribution = ZAP or ZANB
- These models are not readily available in SAS, but NBH is in Mplus
 - Could be fit as a multivariate model in SAS GLIMMIX

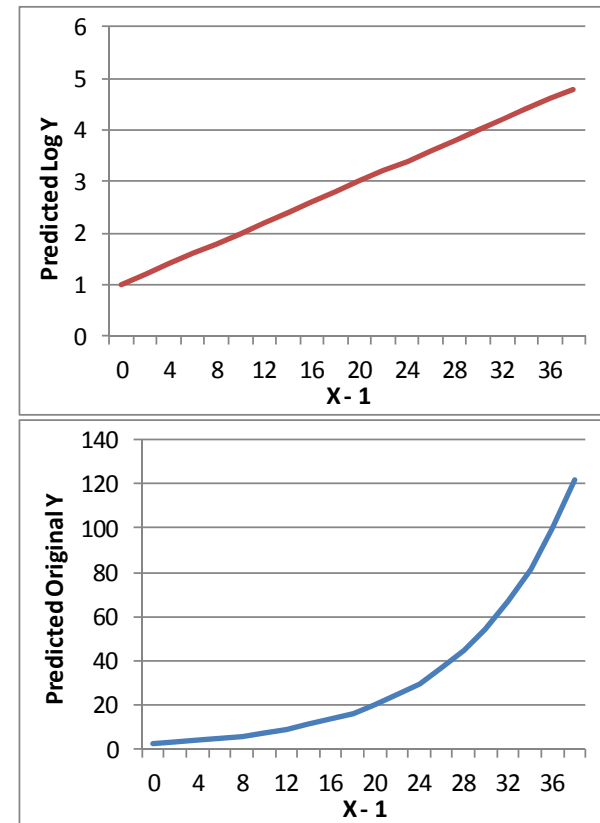
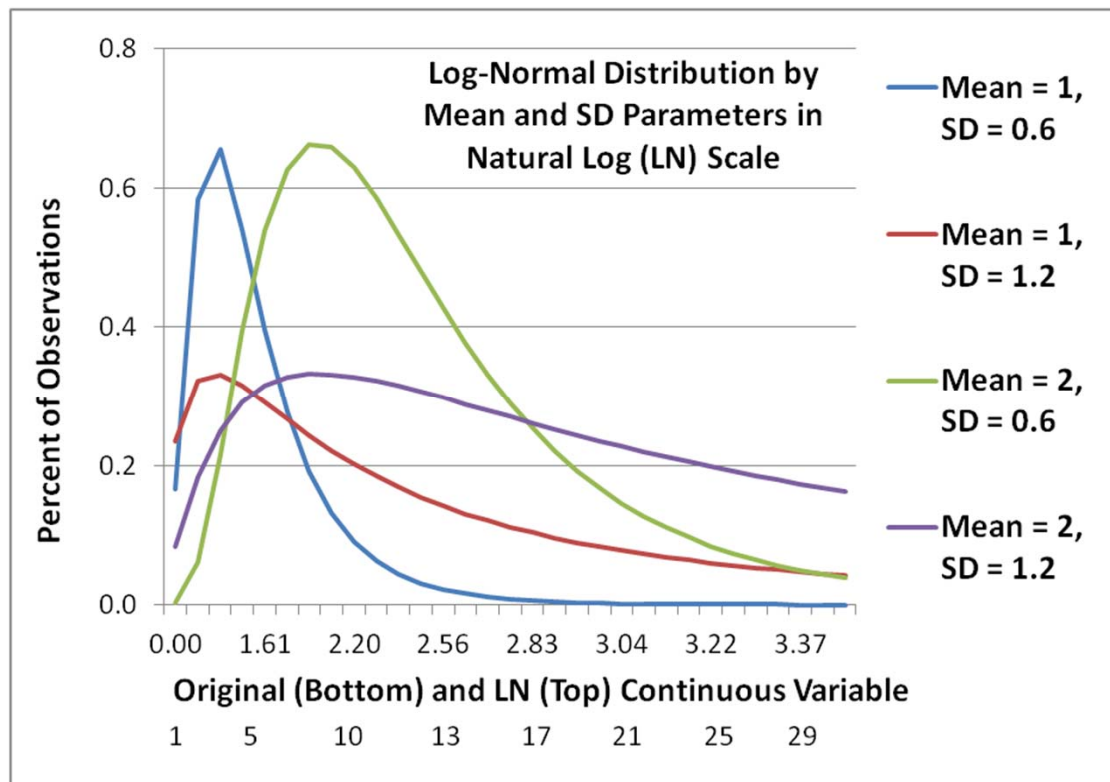
Comparing Models for Count Data

- Whether or not a dispersion parameter is needed can be answered via a likelihood ratio test
 - For the most fair comparison, keep the linear predictor model the same
- Whether or not a zero-inflation model is needed should, in theory, also be answerable via a likelihood ratio test...
 - But people disagree about this
 - Problem? Zero-inflation probability can't be negative, so is bounded at 0
 - Other tests have been proposed (e.g., Vuong test—see SAS macro online)
 - Can always check AIC and BIC (smaller is better)
- In general, models with the same distribution and different links can be compared via AIC and BIC, but one cannot use AIC and BIC to compare across alternative distributions (e.g., normal or not?)
 - Log-Likelihoods are not on the same scale due to using different PDFs
 - Count data can also be modeled using distributions for continuous data...

Models for Continuous Outcomes > 0

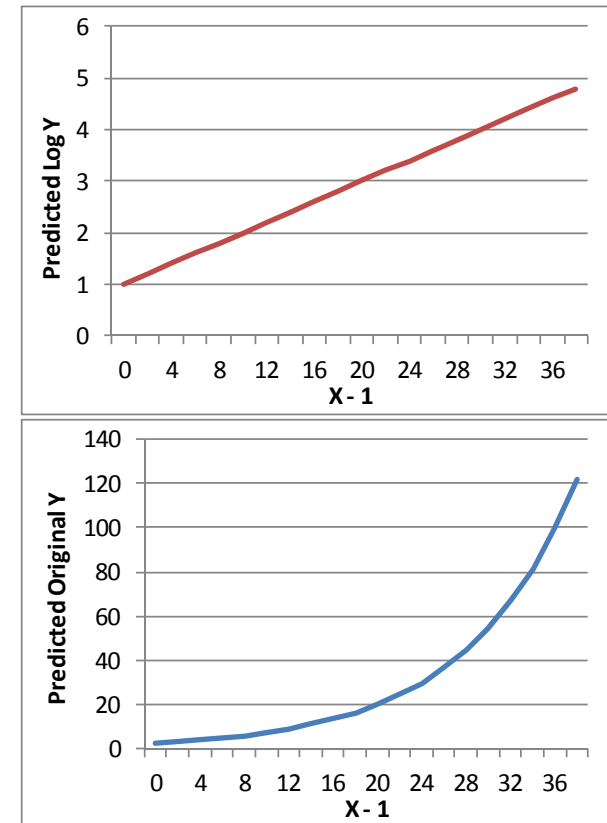
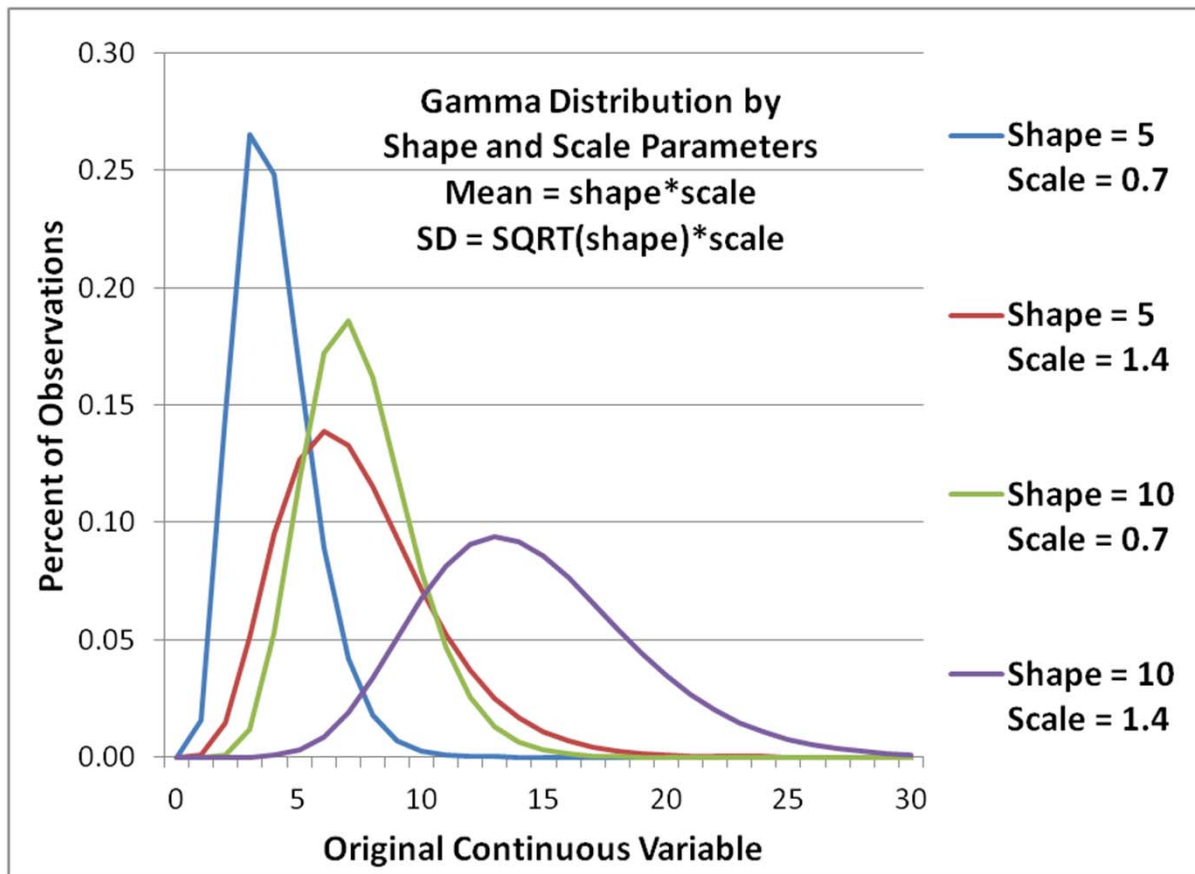
- There are many choices for modeling not-normal *continuous* outcomes (that include positive values only)
 - Most rely on either an identity or log link
 - Will find them in SAS PROC GENMOD and GLIMMIX (see also QLIM)
- GENMOD: DIST= (default link)
 - Gamma (Inverse), Geometric (Log), Inverse Gaussian (Inverse²), Normal (Identity)
- GLIMMIX: DIST= (default link)
 - Beta (Logit), Exponential (Log), Gamma (Log), Geometric (Log), Inverse Gaussian (Inverse²), Normal (Identity), LogNormal (Identity), TCentral (Identity), and BYOBS, which allows multivariate models by which you specify DV-specific models estimated simultaneously (e.g., two-part)
- Many others possible as well—here are just some examples...

Log-Normal Distribution (Link=Identity)



- Model: $\text{Log}(y_i) = \beta_0 + \beta_1 X_i + \beta_2 Z_i + e_i$
 where $e_i \sim \text{LogNormal}(0, \sigma_e^2) \rightarrow$ log of residuals is normal, not residuals
 - Happens to be the same as log-transforming your outcome in this case...
 - The LOG function keeps the predicted values positive, but results in an exponential, not linear prediction of original outcome from slopes
 - GLIMMIX provides "intercept" and "scale=SD" that need to be converted...

Gamma Response Distribution



- Model: $\text{Log}(y_i) = \beta_0 + \beta_1 X_i + \beta_2 Z_i + e_i$
where $e_i \sim \text{Gamma}(0, \sigma_e^2) \rightarrow$ variance is based on shape and scale parameters
 - Default Link is log in GLIMMIX, but inverse in GENMOD
 - Provides "intercept" and "scale=1/scale" that need to be converted...

Two-Part Models for Continuous Outcomes

- A two-part model is an analog to hurdle models for zero-inflated count outcomes (and could be used with count outcomes, too)
 - Explicitly **separates 0 from non-0 values** through a mixture of distributions (Bernoulli + Normal or LogNormal)
 - Creates two submodels to predict “if not 0” and “if not 0, how much?”
 - Easier to think about in terms of prediction (in my opinion)
- Submodel 1: $\text{Logit}(y_i > 0) = \beta_{01} + \beta_{11}X_i + \beta_{21}Z_i$
 - Predict being a 0 using Link = Logit, Distribution = Bernoulli
 - Usually do specify predictors for this part
- Submodel 2: $(y_i | y_i > 0) = \beta_{02} + \beta_{11}X_i + \beta_{21}Z_i$
 - Predict rest of positive amount using Link = Identity, Distribution = Normal or Log-Normal (often rest of distribution is skewed, so log works better)
- Two-part is not readily available in SAS, but is in Mplus
 - Could be fit as a multivariate model in SAS GLIMMIX (I think)
 - Is related to “tobit” models for censored outcomes (for floor/ceiling effects)