

Supporting Materials for ‘The Integration of Norse-Derived Terms in English’

Contents

1	Introduction	1
1.1	Variables	2
2	Load Libraries	3
3	Load data	3
3.1	Examples	4
4	Descriptive statistics	5
5	Total frequency analysis	10
5.1	Simple distance (total frequency)	10
5.2	Feature-based distance (total frequency)	17
5.3	Historical distance (total frequency)	20
5.4	Exploratory analyses	24
6	Source-level frequency analysis	30
6.1	Simple distance	32
6.2	Feature-based distance	37
6.3	Historical distance	42
7	Comparison between texts	51

1 Introduction

The analysis covers the following sources:

- FCPC Peterborough, 1154 (12th C)
- Ormulum South Lincolnshire, ca. 1180 (12th C)
- Havelok, composed in Lincolnshire, but manuscript from West Norfolk 13th century (13th C)
- Genesis and Exodus, composed in East Midlands, manuscript from West Norfolk 13th/14th century (14th C)
- Cursor mundi, c. 1300 (14th C)
- Gawain-poet Cheshire / Staffordshire, ca. 1380 (14th C)
- St Erkenwald, written late fourteenth or the early fifteenth century, manuscript from 1477 (15th C)
- Mannyng South Lincolnshire, ca. 1450 (15th C)
- Wars of Alexander northern England, ca. 1450 (15th C)
- Texts from Lincolnshire, Norfolk and Nottinghamshire from the Corpus of Middle English (1399-1525), (15th C)
- Texts by Richard Rolle (before 1465), (15thC)

The data analysed here is the product of several processing steps. The file `data/SharedIntegrationOfCognatesData.xls` has the original transcribed data. This is cleaned and processed by the script `analysis/analyseTextDistances.py`, which draws some code from `analysis/CLTSFeatureBasedAlignment.py` which is mainly contributed by Johann Mattis List (see <https://calc.hypotheses.org/1962> and <https://gist.github.com/LinguList/7fac44813572f65259c872ef89fa64ad>). The script calculates the distances between pairs of Norse and English forms according to three measures:

- Simple distance from Keller (2023).
- A feature-based distance.
- A historical distance that uses the likelihood of one segment historically replacing another.

Each row in the data represents a comparison between a Norse form and an English form within a given cognate set, including the three measures of distance and frequency of occurrence in each source.

1.1 Variables

The variables in the data (and some that are calculated in the script below) are described as follows:

“Set”: the cognate set the pair of comparisons belong to.

“Class”: the word class the set can appear as (and some binary variables representing the same information).

“NorseLexeme”, “EngLexeme”, “NorseForm”, “EngForm”: The lexemes and full forms for Norse and English terms.

“NorseFormDiagnostic”, “EngFormDiagnostic”: the relevant parts of the form that are diagnostic of the etymology.

“Alignment”, “FeatureAlignment”, “HistoricalAlignment”: multiple sequence alignment of the forms according to the three measures.

“RawDistance”, “NormDistance”, “RawFeatureDistance”, “NormFeatureDistance”, “RawHistoricalDistance”, “NormHistoricalDistance”: Raw and normed distances between

“NFreq...”, “EFreq...”: Frequencies of the Norse and English forms in each source.

“Alliteration”: only true if the source uses alliterative verse and the Norse form and the English form do *not* start with segments in the same alliterative category. That is, it is true if part of the decision about which form to use might be influenced by the need to alliterate. The alliterative texts include: Gawain, St Erkenwold, Wars of Alexander.

“totalNFreq”: The total Norse frequency across all sources.

“totalEFreq”: The total English frequency across all sources.

“NorseProp”: The proportion of total Norse frequency compared to total Norse Frequency + total English Frequency.

“NorseDiagnosticScore”: Whether or not the forms differ according to a specific segments that are characteristically diagnostic of Norse etymology.

2 Load Libraries

```
library(openxlsx)
library(sjPlot)
library(lme4)
library(mgcv)
library(party)
library(ggplot2)
library(phangorn)
library(gridExtra)
library(GGally)
library(MuMIn)
library(brms)
library(ggeffects)
```

3 Load data

Load data created by the python program and convert variables to their proper type:

```
d = read.xlsx("../data/IntegrationDistances.xlsx",1)
# Ignore numerals in set name
d$Set = gsub("[0-9]", "", d$Set)
d$Set = as.factor(d$Set)
d$EngLexeme = factor(d$EngLexeme)
d$ELen = nchar(d$EngForm)
d$NLen = nchar(d$NorseForm)

norseFrequencyColumns = c("NFreqFCPC", "NFreqGawainPoet",
                          "NFreqGenAndEx", "NFreqHavelok",
                          "NFreqMannyng", "NFreqOrmulum",
                          'NFreqWarsAlexander',
                          "NFreqStErkenwald", "NFreqCursorMundi",
                          "NFreqLinc", "NFreqNott", "NFreqNorf", "NFreqRolle")

englishFrequencyColumns = c("EFreqFCPC", "EFreqGawainPoet",
                           "EFreqGenAndEx", "EFreqHavelok",
                           "EFreqMannyng", "EFreqOrmulum",
                           "EFreqWarsAlexander",
                           "EFreqStErkenwald", "EFreqCursorMundi",
                           "EFreqLinc", "EFreqNott", "EFreqNorf", "EFreqRolle")
```

Convert frequencies to numeric type:

```
for(col in c(norseFrequencyColumns, englishFrequencyColumns)){
  d[,col] = as.numeric(d[,col])
}
```

Calculate the total frequency across all sources and the proportion of Norse forms compared to all forms:

```
d$totalNFreq = rowSums(d[,norseFrequencyColumns], na.rm = T)
d$totalEFreq = rowSums(d[,englishFrequencyColumns], na.rm = T)
d$NorseProp = d$totalNFreq / rowSums(d[,c("totalNFreq", "totalEFreq")], na.rm = T)

d = d[!is.na(d$NorseProp),]
```

Scale the distances to lie between 0 and 1:

```
d$NormDistance = as.numeric(d$NormDistance)
d$NormFeatureDistance = as.numeric(d$NormFeatureDistance)
d$NormHistoricalDistance = as.numeric(d$NormHistoricalDistance)
```

```

# Scale distances
normX = function(X){
  (X-min(X))/(max(X)-min(X))
}
d$NormDistance = normX(d$NormDistance)
d$NormFeatureDistance = normX(d$NormFeatureDistance)
d$NormHistoricalDistance = normX(d$NormHistoricalDistance)

d$NormDistance.rank =
  rank(d$NormDistance,ties.method = "max")
d$NormFeatureDistance.rank =
  rank(d$NormFeatureDistance,ties.method = "max")
d$NormHistoricalDistance.rank =
  rank(d$NormHistoricalDistance,ties.method = "max")

```

3.1 Examples

4 Descriptive statistics

Number of English forms: 135

Number of Norse forms: 127

Number of comparisons: 1638

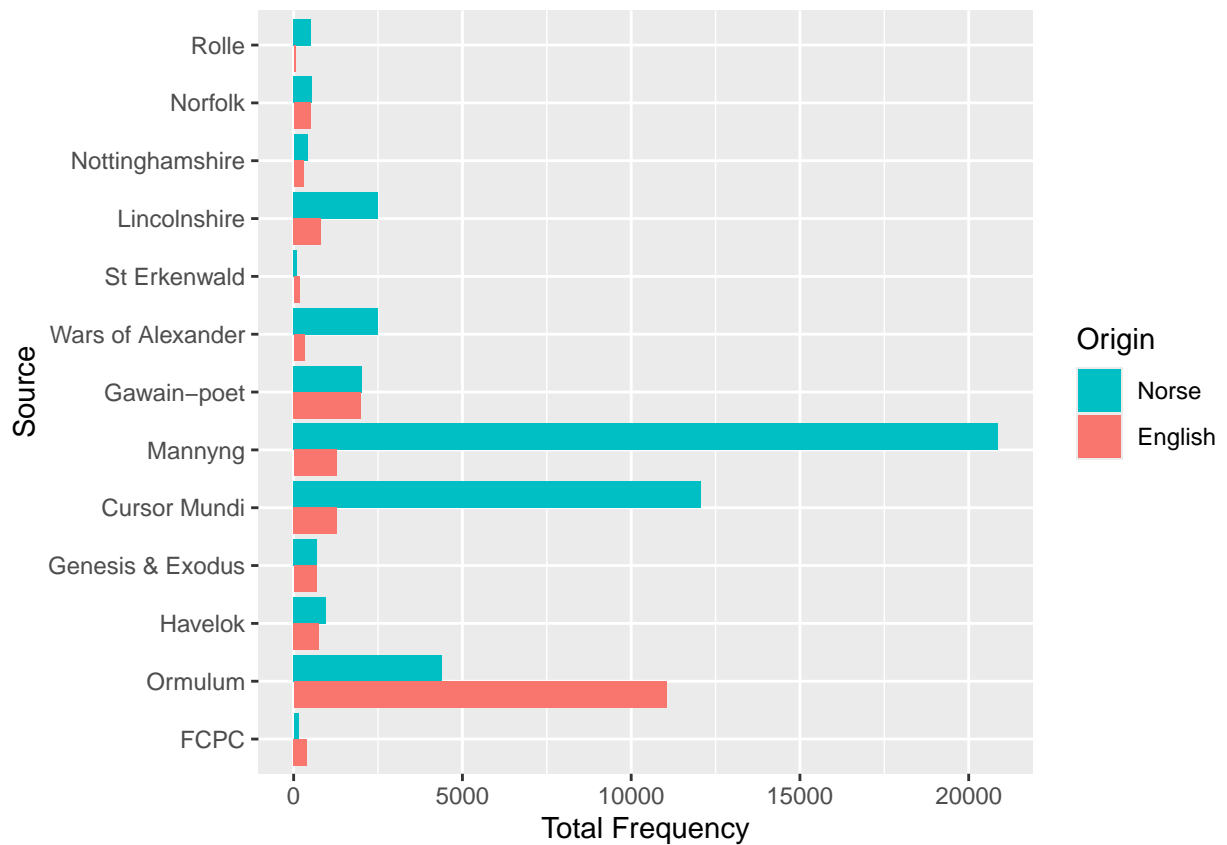
Number of sets: 67

Proportion of Norse and English terms (relative total frequency of each type):

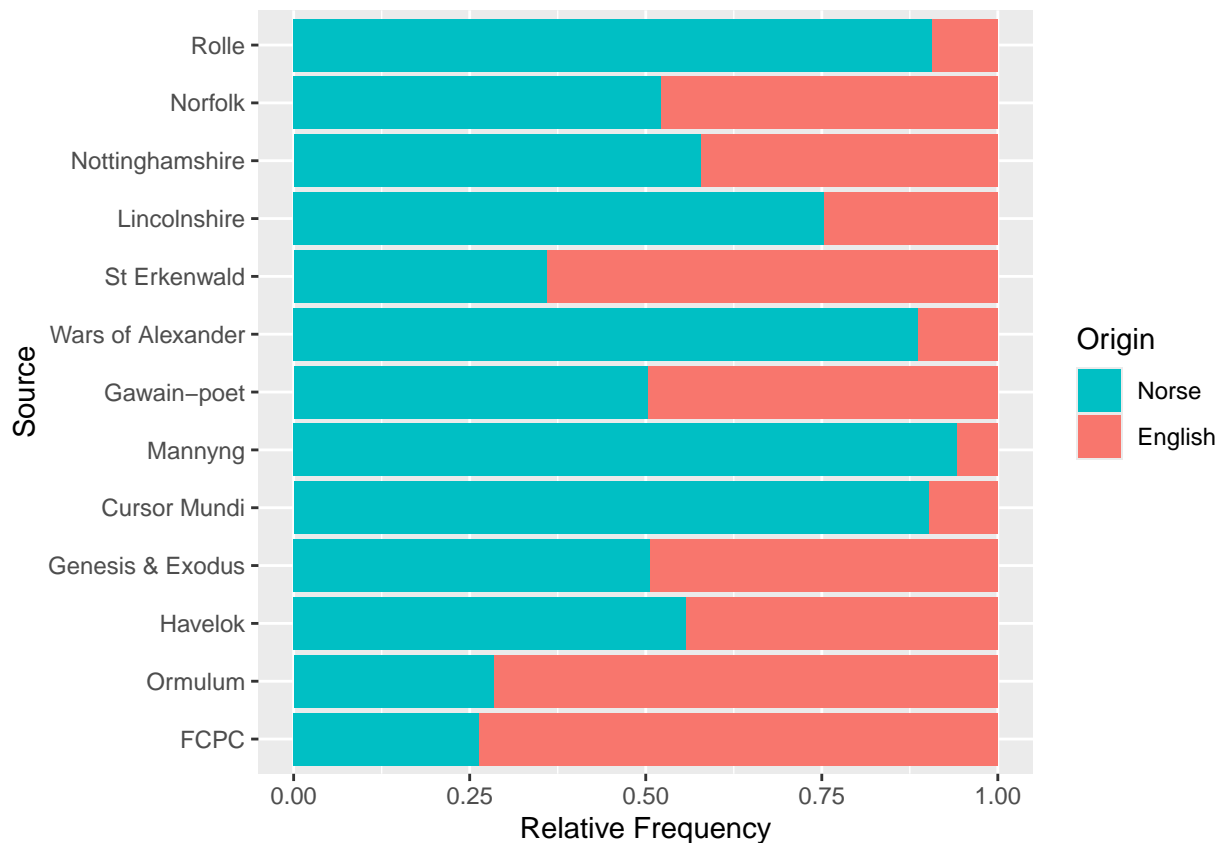
```
dx = rbind(
  data.frame(
    Lang = "Norse",
    Freq = colSums(d[,norseFrequencyColumns],na.rm = T),
    Source = c("FCPC","Gawain-poet","Genesis & Exodus","Havelok",
               "Mannyng", "Ormulum","Wars of Alexander",
               "St Erkenwald", "Cursor Mundi",
               "Lincolnshire","Nottinghamshire",
               "Norfolk", "Rolle")),
  data.frame(
    Lang = "English",
    Freq = colSums(d[,englishFrequencyColumns],na.rm = T),
    Source = c("FCPC","Gawain-poet","Genesis & Exodus","Havelok",
               "Mannyng", "Ormulum","Wars of Alexander",
               "St Erkenwald", "Cursor Mundi",
               "Lincolnshire","Nottinghamshire",
               "Norfolk", "Rolle")))

dx$Source = factor(dx$Source,
  levels=c("FCPC","Ormulum","Havelok",
           "Genesis & Exodus",
           "Cursor Mundi","Mannyng",
           "Gawain-poet", "Wars of Alexander",
           "St Erkenwald",
           "Lincolnshire","Nottinghamshire",
           "Norfolk", "Rolle"))

gRawFreq = ggplot(dx,aes(x=Source,y=Freq,fill=Lang)) +
  geom_bar(stat="identity",position = 'dodge') +
  ylab("Total Frequency") +
  scale_fill_discrete(breaks=c('Norse',"English"),
                      name = "Origin") +
  coord_flip()
gRawFreq
```



```
gPropFreq = ggplot(dx,aes(x=Source,y=Freq,fill=Lang)) +
  geom_bar(stat="identity",position = 'fill') +
  ylab("Relative Frequency") +
  scale_fill_discrete(breaks=c('Norse',"English"),
                      name = "Origin")+
  coord_flip()
gPropFreq
```



```
pdf("../results/GRawFreq.pdf",width=5,height=3.5)
gRawFreq
dev.off()
```

```
## pdf
## 2
```

```
pdf("../results/GPropFreq.pdf",width=5,height=3.5)
gPropFreq
dev.off()
```

```
## pdf
## 2
```

```
pdf("../results/GRawAndPropFreq.pdf",width=6,height=3)
grid.arrange(gRawFreq+
  theme(legend.position = "none"),
  gPropFreq +
  theme(axis.text.y = element_blank(),
        axis.title.y = element_blank()),
  ncol=2)
dev.off()
```

```
## pdf
## 2
```

Correlations between distance measures:

```
cor.test(d$NormDistance,d$NormFeatureDistance)
```

```
##
## Pearson's product-moment correlation
##
## data: d$NormDistance and d$NormFeatureDistance
## t = 94.868, df = 1636, p-value < 2.2e-16
```

```
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
## 0.9120847 0.9270142
## sample estimates:
## cor
## 0.919882
```

```
cor.test(d$NormDistance,d$NormHistoricalDistance)
```

```
##
## Pearson's product-moment correlation
##
## data: d$NormDistance and d$NormHistoricalDistance
## t = 75.069, df = 1636, p-value < 2.2e-16
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
## 0.8689625 0.8907967
## sample estimates:
## cor
## 0.8803451
```

```
cor.test(d$NormFeatureDistance,d$NormHistoricalDistance)
```

```
##
## Pearson's product-moment correlation
##
## data: d$NormFeatureDistance and d$NormHistoricalDistance
## t = 75.886, df = 1636, p-value < 2.2e-16
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
## 0.8712788 0.8927493
## sample estimates:
## cor
## 0.8824729
```

Function for extracting stats:

```
getStatReport = function(m0,m1,m2,m3,finalModel,outFile){
  modelComparison = as.data.frame(anova(m0,m1,m2,m3))
  modelComparison$pChi = round(modelComparison$`Pr(>Chisq)`,3)
  modelComparison$pChi[modelComparison$pChi==0] = "< 0.001"
  modelComparison$lldiff = NA
  modelComparison$lldiff[2:4] = diff(modelComparison$logLik)
  modelComparison = modelComparison[2:4,]
  modelComparison$Term = c("Linear","Quadratic","Cubic")

  coef = as.data.frame(summary(finalModel)$coefficients)
  rx = which(grepl("Norm",rownames(coef)))
  coef[rx,"Estimate"] = round(coef[rx,"Estimate"],3)
  coef[rx,"z value"] = round(coef[rx,"z value"],3)
  coef[, "Pr(>|z|)"] = round(coef[, "Pr(>|z|)"],3)
  coef[, "Pr(>|z|)"][coef[, "Pr(>|z|)"] ==0] = "< 0.001"

  modelComparison$Estimate = coef[rx,"Estimate"]
  modelComparison$z = coef[rx,"z value"]
  modelComparison$p = coef[rx,"Pr(>|z|)"]

  modelComparison = modelComparison[,
    c("Term","BIC","lldiff",
      "Chisq","Df","pChi",
      "Estimate","z","p")]
}
```



```

mcx = paste(paste0(c("Linear: ", "Quadratic: ", "Cubic: "),
  "beta = ", coef[rx, "Estimate"],
  ", z = ", coef[rx, "z value"],
  ", Wald p = ", coef[rx, "Pr(>|z|)"],
  ", LLDiff = ", round(modelComparison$lldiff, 1),
  ", df = ", modelComparison$Df,
  ", p = ", modelComparison$pChi), collapse="; ")
mcx = gsub("= <", "<", mcx)
cat(mcx, file=outFile)

# mx = rbind(c("", "Model Comparison", "", "", "", "", "Model Estimate", "", ""),
#           names(modelComparison),
#           modelComparison)

modelComparison$BIC = round(modelComparison$BIC, 1)
modelComparison$lldiff = round(modelComparison$lldiff, 1)
modelComparison$Chisq = round(modelComparison$Chisq, 1)
modelComparison$Estimate = round(modelComparison$Estimate, 2)

write.csv(modelComparison, file=gsup("\\.txt", ".csv", outFile))

return(mcx)
}

```

5 Total frequency analysis

The analyses below predict the total frequency across all sources, using the orthographic form as the basis for observations.

5.1 Simple distance (total frequency)

We use nested model comparison to evaluate whether higher-order variables should be added (i.e. how non-linear the relationships is).

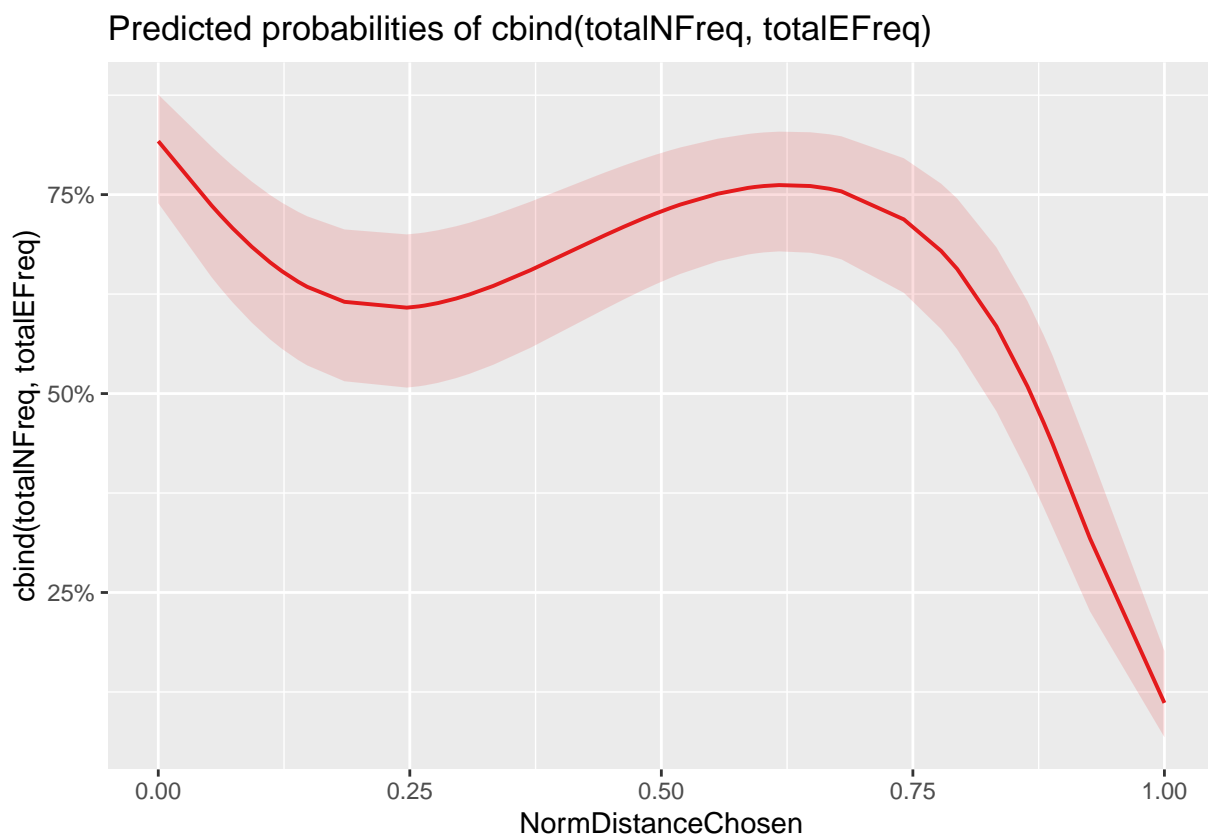
```
d$NormDistanceChosen = d$NormDistance
# Baseline model
m0 = glmer(cbind(totalNFreq,totalEFreq) ~
           1 + (1|Set) + (1|EngLexeme),
           data = d, family = "binomial")
# Add the distance measure
m1 = update(m0,~.+NormDistanceChosen)
# Add quadratic term
m2 = update(m1,~.+I(NormDistanceChosen^2))
# Add cubic term
m3 = update(m2,~.+I(NormDistanceChosen^3))
# Compare fit of models
anova(m0,m1,m2,m3)

## Data: d
## Models:
## m0: cbind(totalNFreq, totalEFreq) ~ 1 + (1 | Set) + (1 | EngLexeme)
## m1: cbind(totalNFreq, totalEFreq) ~ (1 | Set) + (1 | EngLexeme) + NormDistanceChosen
## m2: cbind(totalNFreq, totalEFreq) ~ (1 | Set) + (1 | EngLexeme) + NormDistanceChosen + I(NormDist
## m3: cbind(totalNFreq, totalEFreq) ~ (1 | Set) + (1 | EngLexeme) + NormDistanceChosen + I(NormDist
##      npar    AIC    BIC logLik deviance   Chisq Df Pr(>Chisq)
## m0      3 32262 32278 -16128    32256
## m1      4 32264 32286 -16128    32256  0.0038  1    0.9511
## m2      5 32207 32234 -16098    32197 59.1556  1 1.457e-14 ***
## m3      6 31975 32007 -15981    31963 234.1472  1 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Summary of final model:
summary(m3)

## Generalized linear mixed model fit by maximum likelihood (Laplace
## Approximation) [glmerMod]
## Family: binomial ( logit )
## Formula: cbind(totalNFreq, totalEFreq) ~ (1 | Set) + (1 | EngLexeme) +
##          NormDistanceChosen + I(NormDistanceChosen^2) + I(NormDistanceChosen^3)
## Data: d
##
##      AIC      BIC   logLik deviance df.resid
## 31974.6 32007.0 -15981.3 31962.6    1632
##
## Scaled residuals:
##      Min       1Q   Median       3Q      Max
## -13.6620  -2.5875  -0.6251   0.7106  24.4508
##
## Random effects:
## Groups   Name                Variance Std.Dev.
## EngLexeme (Intercept) 1.997      1.413
## Set      (Intercept) 1.462      1.209
```

```
## Number of obs: 1638, groups: EngLexeme, 135; Set, 67
##
## Fixed effects:
##               Estimate Std. Error z value Pr(>|z|)
## (Intercept)      1.4975    0.2323   6.445 1.15e-10 ***
## NormDistanceChosen -10.4393    0.9350 -11.165 < 2e-16 ***
## I(NormDistanceChosen^2) 30.8294    2.1940  14.052 < 2e-16 ***
## I(NormDistanceChosen^3) -23.9658    1.5108 -15.862 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Correlation of Fixed Effects:
##               (Intr) NrmDsC I(NDC^2
## NrmDstncChs -0.432
## I(NrmDsC^2)  0.384 -0.966
## I(NrmDsC^3) -0.342  0.904 -0.981
pSimpleTotal = plot_model(m3,'eff',
  terms="NormDistanceChosen[all]")
pSimpleTotal
```



Marginal represents the variance explained by the fixed effects. Conditional represents the variance explained by the entire model.

```
r.squaredGLMM(m3)
```

```
## Warning: 'r.squaredGLMM' now calculates a revised statistic. See the help page.
## Warning: the null model is correct only if all variables used by the original
## model remain unchanged.
##               R2m      R2c
## theoretical 0.04430528 0.9783727
## delta      0.04404589 0.9726447
```

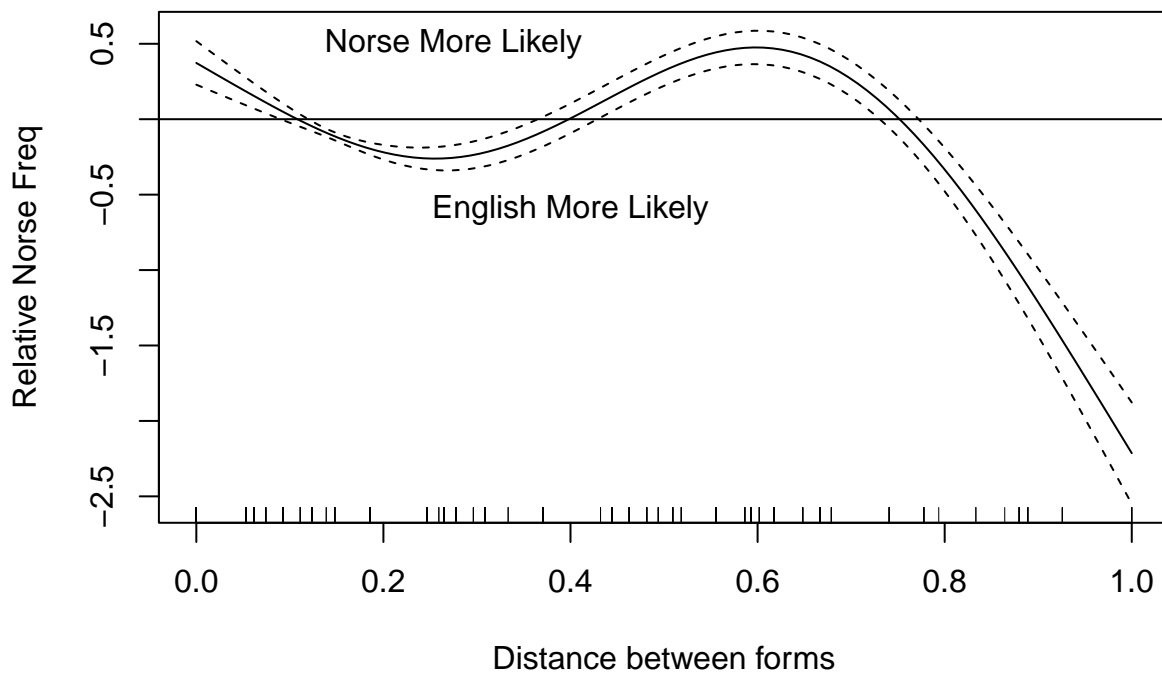
Statistics: (Linear: beta = -10.439, z = -11.165, Wald p < 0.001, LLDiff = 0, df = 1, p = 0.951; Quadratic: beta = 30.829, z = 14.052, Wald p < 0.001, LLDiff = 29.6, df = 1, p < 0.001; Cubic: beta = -23.966, z = -15.862, Wald p < 0.001, LLDiff = 117.1, df = 1, p < 0.001)

Same analysis using GAM with binomial distribution:

```
mGAM = gam(cbind(totalNFreq,totaleFreq) ~
            s(NormDistanceChosen,k=4) +
            s(Set,bs="re") +
            s(EngLexeme,bs="re"),
            data = d, family = "binomial")
summary(mGAM)

##
## Family: binomial
## Link function: logit
##
## Formula:
## cbind(totalNFreq, totaleFreq) ~ s(NormDistanceChosen, k = 4) +
##      s(Set, bs = "re") + s(EngLexeme, bs = "re")
##
## Parametric coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  0.6651      0.2292   2.902  0.00371 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Approximate significance of smooth terms:
##              edf Ref.df   Chi.sq p-value
## s(NormDistanceChosen)  2.998     3    250.5 <2e-16 ***
## s(Set)                44.361    66 1409276.0  0.361
## s(EngLexeme)          81.753   134 355114.5  0.545
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## R-sq.(adj) =  0.438   Deviance explained = 45.9%
## UBRE = 16.108   Scale est. = 1           n = 1638

plot.gam(mGAM,
          ylab="Relative Norse Freq",
          xlab="Distance between forms",
          select = 1)
abline(h=0)
text(0.275,0.5,"Norse More Likely")
text(0.4,-0.6,"English More Likely")
```



The raw frequency of Norse forms is negatively correlated with distance:

```
mNorseFreq = glmer(totalNFreq ~ NormDistance +
  (1|Set) + (1|NorseLexeme),
  data = d, family = "poisson")
```

The distribution of Norse proportions is bimodal, with most terms having either low or high proportions. This isn't ideal for the binomial models above, so we also test the same model using a beta binomial distribution, which can fit bimodal binomial distributions. The result is very similar to the others:

```
d$totalAllFreq = d$totalNFreq + d$totalEFreq
mBinomBeta = brm(totalNFreq | trials(totalAllFreq) ~
  1 + NormDistanceChosen +
  I(NormDistanceChosen^2) +
  I(NormDistanceChosen^3) +
  (1|Set) + (1|EngLexeme),
  data = d, family = "beta_binomial")
```

```
## Compiling Stan program...
```

```
## Start sampling
```

```
##
```

```
## SAMPLING FOR MODEL '7404dd55f235b75d82269cea3d068515' NOW (CHAIN 1).
```

```
## Chain 1:
```

```
## Chain 1: Gradient evaluation took 0.000587 seconds
```

```
## Chain 1: 1000 transitions using 10 leapfrog steps per transition would take 5.87 seconds.
```

```
## Chain 1: Adjust your expectations accordingly!
```

```
## Chain 1:
```

```
## Chain 1:
```

```
## Chain 1: Iteration: 1 / 2000 [ 0%] (Warmup)
```

```
## Chain 1: Iteration: 200 / 2000 [ 10%] (Warmup)
```

```
## Chain 1: Iteration: 400 / 2000 [ 20%] (Warmup)
```

```
## Chain 1: Iteration: 600 / 2000 [ 30%] (Warmup)
```

```
## Chain 1: Iteration: 800 / 2000 [ 40%] (Warmup)
```

```
## Chain 1: Iteration: 1000 / 2000 [ 50%] (Warmup)
```

```
## Chain 1: Iteration: 1001 / 2000 [ 50%] (Sampling)
```

```
## Chain 1: Iteration: 1200 / 2000 [ 60%] (Sampling)
```

```
## Chain 1: Iteration: 1400 / 2000 [ 70%] (Sampling)
```

```

## Chain 1: Iteration: 1600 / 2000 [ 80%] (Sampling)
## Chain 1: Iteration: 1800 / 2000 [ 90%] (Sampling)
## Chain 1: Iteration: 2000 / 2000 [100%] (Sampling)
## Chain 1:
## Chain 1: Elapsed Time: 78.4587 seconds (Warm-up)
## Chain 1: 118.942 seconds (Sampling)
## Chain 1: 197.401 seconds (Total)
## Chain 1:
##
## SAMPLING FOR MODEL '7404dd55f235b75d82269cea3d068515' NOW (CHAIN 2).
## Chain 2:
## Chain 2: Gradient evaluation took 0.000507 seconds
## Chain 2: 1000 transitions using 10 leapfrog steps per transition would take 5.07 seconds.
## Chain 2: Adjust your expectations accordingly!
## Chain 2:
## Chain 2:
## Chain 2: Iteration: 1 / 2000 [ 0%] (Warmup)
## Chain 2: Iteration: 200 / 2000 [ 10%] (Warmup)
## Chain 2: Iteration: 400 / 2000 [ 20%] (Warmup)
## Chain 2: Iteration: 600 / 2000 [ 30%] (Warmup)
## Chain 2: Iteration: 800 / 2000 [ 40%] (Warmup)
## Chain 2: Iteration: 1000 / 2000 [ 50%] (Warmup)
## Chain 2: Iteration: 1001 / 2000 [ 50%] (Sampling)
## Chain 2: Iteration: 1200 / 2000 [ 60%] (Sampling)
## Chain 2: Iteration: 1400 / 2000 [ 70%] (Sampling)
## Chain 2: Iteration: 1600 / 2000 [ 80%] (Sampling)
## Chain 2: Iteration: 1800 / 2000 [ 90%] (Sampling)
## Chain 2: Iteration: 2000 / 2000 [100%] (Sampling)
## Chain 2:
## Chain 2: Elapsed Time: 84.889 seconds (Warm-up)
## Chain 2: 68.5895 seconds (Sampling)
## Chain 2: 153.478 seconds (Total)
## Chain 2:
##
## SAMPLING FOR MODEL '7404dd55f235b75d82269cea3d068515' NOW (CHAIN 3).
## Chain 3:
## Chain 3: Gradient evaluation took 0.000479 seconds
## Chain 3: 1000 transitions using 10 leapfrog steps per transition would take 4.79 seconds.
## Chain 3: Adjust your expectations accordingly!
## Chain 3:
## Chain 3:
## Chain 3: Iteration: 1 / 2000 [ 0%] (Warmup)
## Chain 3: Iteration: 200 / 2000 [ 10%] (Warmup)
## Chain 3: Iteration: 400 / 2000 [ 20%] (Warmup)
## Chain 3: Iteration: 600 / 2000 [ 30%] (Warmup)
## Chain 3: Iteration: 800 / 2000 [ 40%] (Warmup)
## Chain 3: Iteration: 1000 / 2000 [ 50%] (Warmup)
## Chain 3: Iteration: 1001 / 2000 [ 50%] (Sampling)
## Chain 3: Iteration: 1200 / 2000 [ 60%] (Sampling)
## Chain 3: Iteration: 1400 / 2000 [ 70%] (Sampling)
## Chain 3: Iteration: 1600 / 2000 [ 80%] (Sampling)
## Chain 3: Iteration: 1800 / 2000 [ 90%] (Sampling)
## Chain 3: Iteration: 2000 / 2000 [100%] (Sampling)
## Chain 3:
## Chain 3: Elapsed Time: 75.5781 seconds (Warm-up)
## Chain 3: 78.028 seconds (Sampling)
## Chain 3: 153.606 seconds (Total)
## Chain 3:

```

```
##
## SAMPLING FOR MODEL '7404dd55f235b75d82269cea3d068515' NOW (CHAIN 4).
## Chain 4:
## Chain 4: Gradient evaluation took 0.000488 seconds
## Chain 4: 1000 transitions using 10 leapfrog steps per transition would take 4.88 seconds.
## Chain 4: Adjust your expectations accordingly!
## Chain 4:
## Chain 4:
## Chain 4: Iteration:    1 / 2000 [  0%] (Warmup)
## Chain 4: Iteration:   200 / 2000 [ 10%] (Warmup)
## Chain 4: Iteration:   400 / 2000 [ 20%] (Warmup)
## Chain 4: Iteration:   600 / 2000 [ 30%] (Warmup)
## Chain 4: Iteration:   800 / 2000 [ 40%] (Warmup)
## Chain 4: Iteration:  1000 / 2000 [ 50%] (Warmup)
## Chain 4: Iteration: 1001 / 2000 [ 50%] (Sampling)
## Chain 4: Iteration: 1200 / 2000 [ 60%] (Sampling)
## Chain 4: Iteration: 1400 / 2000 [ 70%] (Sampling)
## Chain 4: Iteration: 1600 / 2000 [ 80%] (Sampling)
## Chain 4: Iteration: 1800 / 2000 [ 90%] (Sampling)
## Chain 4: Iteration: 2000 / 2000 [100%] (Sampling)
## Chain 4:
## Chain 4: Elapsed Time: 75.2296 seconds (Warm-up)
## Chain 4:           86.3069 seconds (Sampling)
## Chain 4:           161.537 seconds (Total)
## Chain 4:
summary(mBinomBeta)

## Family: beta_binomial
## Links: mu = logit; phi = identity
## Formula: totalNFreq | trials(totalAllFreq) ~ 1 + NormDistanceChosen + I(NormDistanceChosen^2) + I
## Data: d (Number of observations: 1743)
## Draws: 4 chains, each with iter = 2000; warmup = 1000; thin = 1;
## total post-warmup draws = 4000
##
## Group-Level Effects:
## ~EngLexeme (Number of levels: 135)
##           Estimate Est.Error l-95% CI u-95% CI Rhat Bulk_ESS Tail_ESS
## sd(Intercept)    0.64     0.08     0.50     0.80 1.00      843      1711
##
## ~Set (Number of levels: 67)
##           Estimate Est.Error l-95% CI u-95% CI Rhat Bulk_ESS Tail_ESS
## sd(Intercept)    0.35     0.13     0.08     0.59 1.01      462      429
##
## Population-Level Effects:
##           Estimate Est.Error l-95% CI u-95% CI Rhat Bulk_ESS
## Intercept           0.57     0.30    -0.02     1.15 1.00      1161
## NormDistanceChosen   -3.84     2.51    -8.75     1.09 1.00      1124
## INormDistanceChosenE2 11.07     5.82    -0.20    22.46 1.00      1171
## INormDistanceChosenE3 -8.83     3.92   -16.62    -1.25 1.00      1251
## Tail_ESS
## Intercept           1938
## NormDistanceChosen   1659
## INormDistanceChosenE2 1678
## INormDistanceChosenE3 1645
##
## Family Specific Parameters:
##           Estimate Est.Error l-95% CI u-95% CI Rhat Bulk_ESS Tail_ESS
## phi           1.89     0.07     1.75     2.03 1.00      4973      2852
```

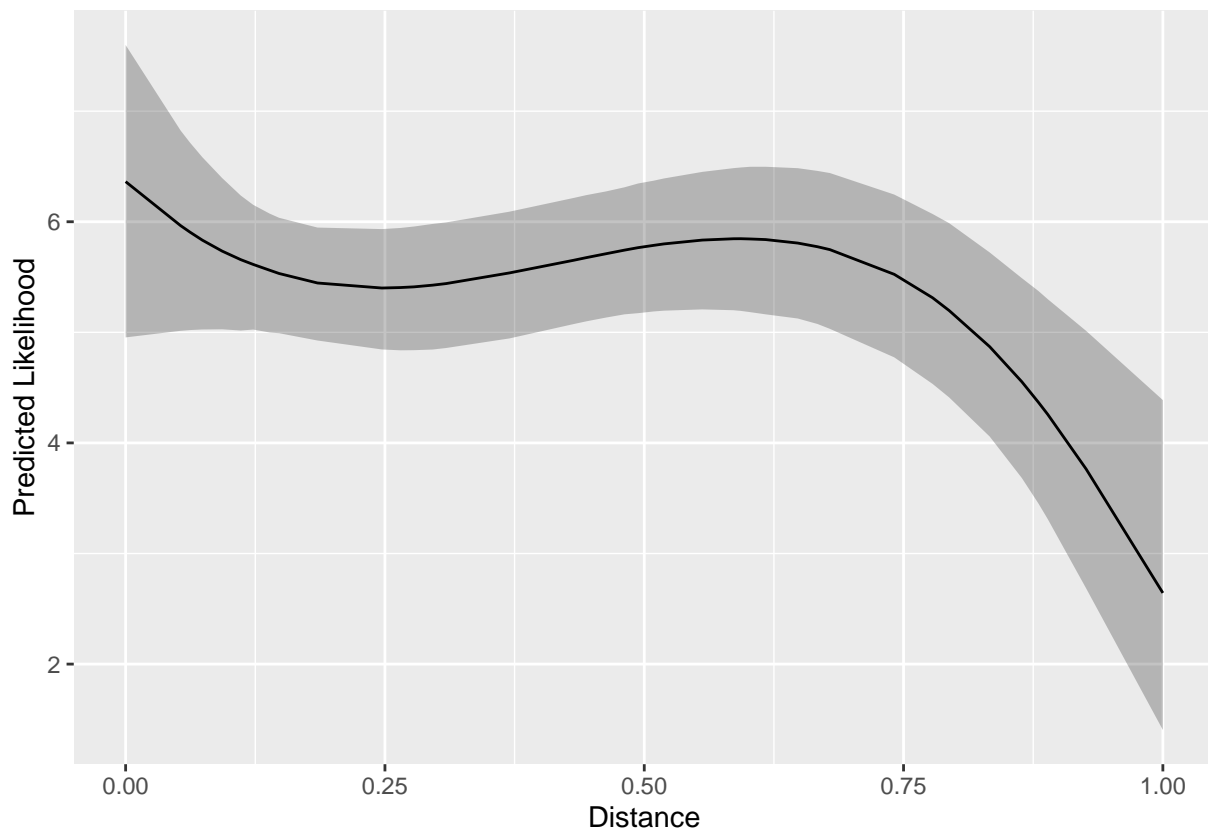
```
##
## Draws were sampled using sampling(NUTS). For each parameter, Bulk_ESS
## and Tail_ESS are effective sample size measures, and Rhat is the potential
## scale reduction factor on split chains (at convergence, Rhat = 1).

pred = ggpredict(mBinomBeta, terms="NormDistanceChosen[all]")

## Note: uncertainty of error terms are not taken into account. Consider
##   setting `interval` to "prediction". This will call `posterior_predict()`
##   instead of `posterior_epred()`.

## Warning in checkMatrixPackageVersion(): Package version inconsistency detected.
## TMB was built with Matrix version 1.6.1
## Current Matrix version is 1.6.5
## Please re-install 'TMB' from source using install.packages('TMB', type = 'source') or ask CRAN for

ggplot(pred, aes(x=x, y=predicted, ymin=conf.low, ymax=conf.high)) +
  geom_ribbon(alpha=0.3) +
  geom_line() +
  xlab("Distance") +
  ylab("Predicted Likelihood")
```



5.2 Feature-based distance (total frequency)

We use nested model comparison to evaluate whether higher-order variables should be added (i.e. how non-linear the relationships is).

```
d$NormDistanceChosen = d$NormFeatureDistance
# Baseline model
m0 = glmer(cbind(totalNFreq,totalEFreq) ~
            1 + (1|Set) + (1|EngLexeme),
            data = d, family = "binomial")
# Add the distance measure
m1 = update(m0,~.+NormDistanceChosen)
# Add quadratic term
m2 = update(m1,~.+I(NormDistanceChosen^2))
# Add cubic term
m3 = update(m2,~.+I(NormDistanceChosen^3))
# Compare fit of models
anova(m0,m1,m2,m3)

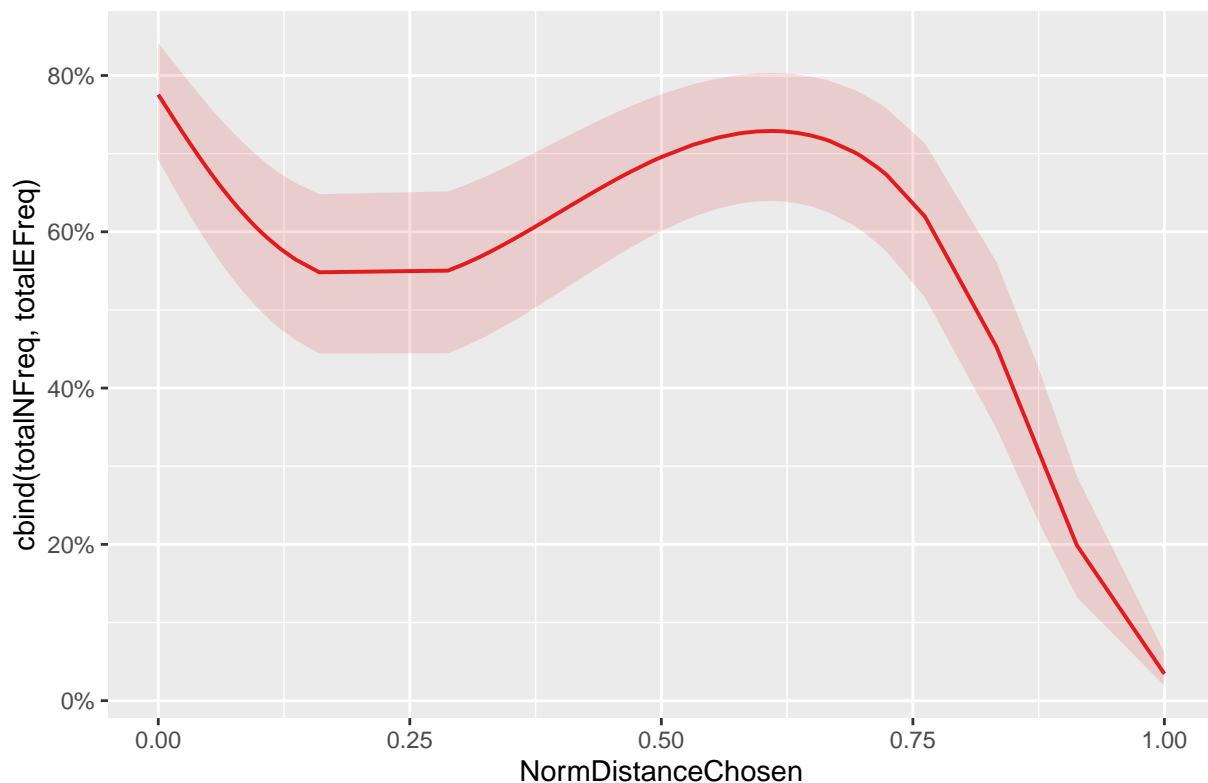
## Data: d
## Models:
## m0: cbind(totalNFreq, totalEFreq) ~ 1 + (1 | Set) + (1 | EngLexeme)
## m1: cbind(totalNFreq, totalEFreq) ~ (1 | Set) + (1 | EngLexeme) + NormDistanceChosen
## m2: cbind(totalNFreq, totalEFreq) ~ (1 | Set) + (1 | EngLexeme) + NormDistanceChosen + I(NormDist
## m3: cbind(totalNFreq, totalEFreq) ~ (1 | Set) + (1 | EngLexeme) + NormDistanceChosen + I(NormDist
##      npar   AIC    BIC logLik deviance   Chisq Df Pr(>Chisq)
## m0      3 34838 34854 -17416    34832
## m1      4 34839 34861 -17416    34831   0.4413  1    0.5065
## m2      5 34766 34793 -17378    34756  75.3682  1    <2e-16 ***
## m3      6 34506 34539 -17247    34494 262.1543  1    <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Summary of final model:
summary(m3)

## Generalized linear mixed model fit by maximum likelihood (Laplace
## Approximation) [glmerMod]
## Family: binomial ( logit )
## Formula: cbind(totalNFreq, totalEFreq) ~ (1 | Set) + (1 | EngLexeme) +
##          NormDistanceChosen + I(NormDistanceChosen^2) + I(NormDistanceChosen^3)
## Data: d
##
##           AIC          BIC    logLik deviance df.resid
## 34505.8 34538.6 -17246.9 34493.8      1737
##
## Scaled residuals:
##      Min       1Q   Median       3Q      Max
## -15.0151  -2.4803  -0.5505   0.7858  24.4498
##
## Random effects:
##  Groups      Name              Variance Std.Dev.
## EngLexeme (Intercept) 1.937      1.392
## Set          (Intercept) 1.550      1.245
## Number of obs: 1743, groups: EngLexeme, 135; Set, 67
##
## Fixed effects:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)      1.2395     0.2188   5.664 1.48e-08 ***
## NormDistanceChosen -11.5412     1.0010 -11.529 < 2e-16 ***
```

```
## I(NormDistanceChosen^2) 35.8979      2.6111  13.748 < 2e-16 ***
## I(NormDistanceChosen^3) -28.9331      1.8700 -15.473 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Correlation of Fixed Effects:
##              (Intr) NrmDsC I(NDC^2
## NrmDstncChs -0.286
## I(NrmDsC^2)  0.266 -0.980
## I(NrmDsC^3) -0.247  0.935 -0.985
pFeatureTotal = plot_model(m3,'eff',
  terms="NormDistanceChosen[all]")
pFeatureTotal
```

Predicted probabilities of cbind(totalNFreq, totalEFreq)



Statistics: (Linear: beta = -11.541, z = -11.529, Wald p < 0.001, LLDiff = 0.2, df = 1, p = 0.506; Quadratic: beta = 35.898, z = 13.748, Wald p < 0.001, LLDiff = 37.7, df = 1, p < 0.001; Cubic: beta = -28.933, z = -15.473, Wald p < 0.001, LLDiff = 131.1, df = 1, p < 0.001)

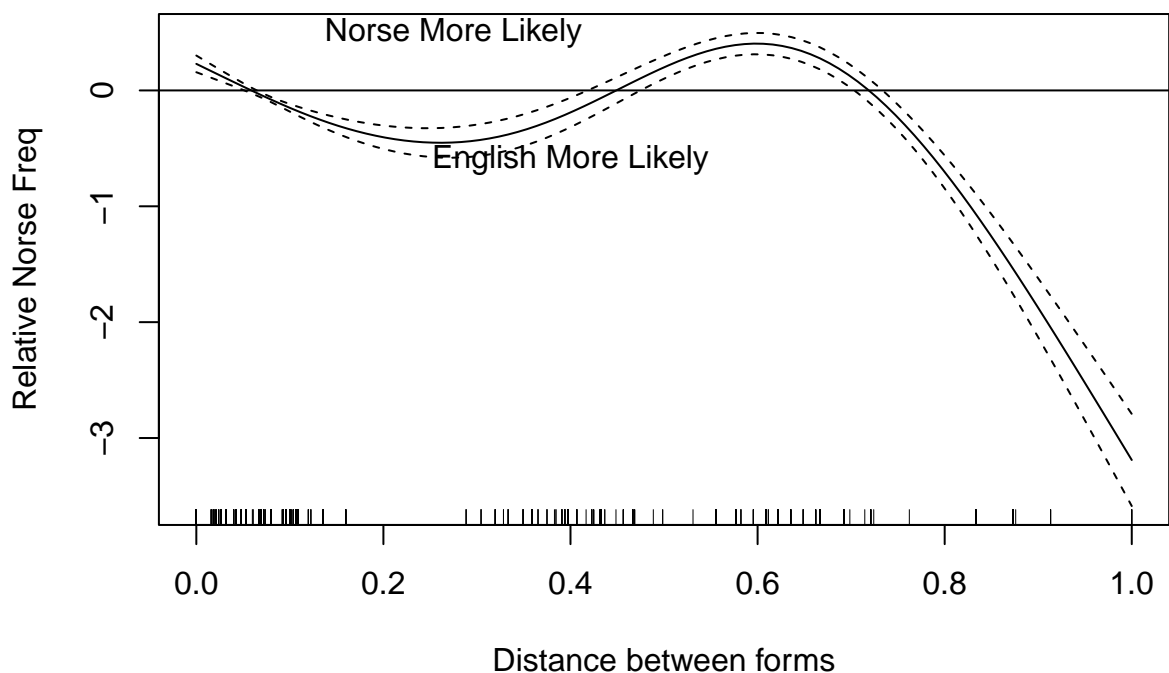
Same analysis using GAM with binomial distribution:

```
mGAM = gam(cbind(totalNFreq,totalEFreq) ~
  s(NormDistanceChosen,k=4) +
  s(Set,bs="re") +
  s(EngLexeme,bs="re"),
  data = d, family = "binomial")
summary(mGAM)
```

```
##
## Family: binomial
## Link function: logit
##
## Formula:
## cbind(totalNFreq, totalEFreq) ~ s(NormDistanceChosen, k = 4) +
```

```
##      s(Set, bs = "re") + s(EngLexeme, bs = "re")
##
## Parametric coefficients:
##           Estimate Std. Error z value Pr(>|z|)
## (Intercept)  0.6475      0.2297   2.819  0.00482 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Approximate significance of smooth terms:
##           edf Ref.df   Chi.sq p-value
## s(NormDistanceChosen)  2.998     3    293.4 <2e-16 ***
## s(Set)                  44.557    66 1317309.0  0.371
## s(EngLexeme)            81.945   134  475900.5  0.503
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## R-sq.(adj) =  0.415   Deviance explained = 43.6%
## UBRE = 16.464   Scale est. = 1         n = 1743
```

```
plot.gam(mGAM,
  ylab="Relative Norse Freq",
  xlab="Distance between forms",
  select = 1)
abline(h=0)
text(0.275,0.5,"Norse More Likely")
text(0.4,-0.6,"English More Likely")
```



5.3 Historical distance (total frequency)

We use nested model comparison to evaluate whether higher-order variables should be added (i.e. how non-linear the relationships is).

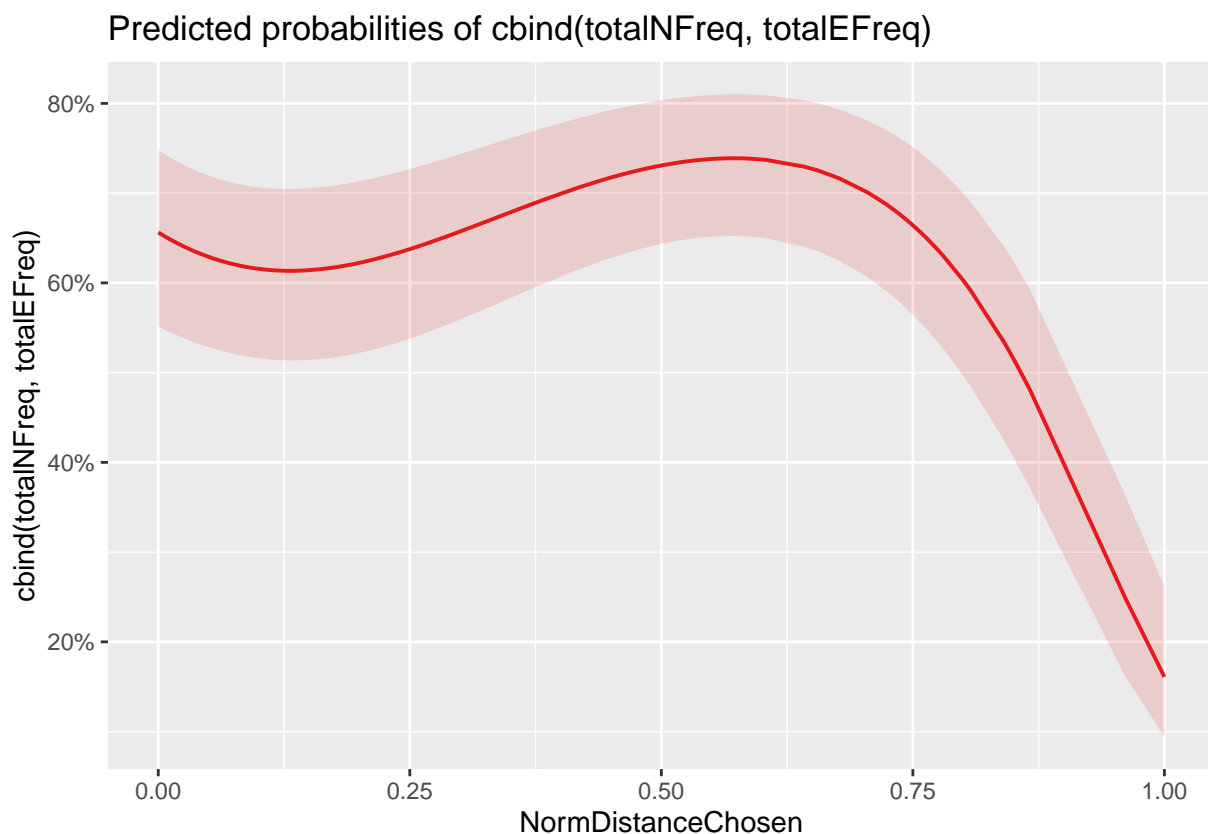
```
d$NormDistanceChosen = d$NormHistoricalDistance
# Baseline model
m0 = glmer(cbind(totalNFreq, totalEFreq) ~
            1 + (1|Set) + (1|EngLexeme),
            data = d, family = "binomial")
# Add the distance measure
m1 = update(m0, ~.+NormDistanceChosen)
# Add quadratic term
m2 = update(m1, ~.+I(NormDistanceChosen^2))
# Add cubic term
m3 = update(m2, ~.+I(NormDistanceChosen^3))
# Compare fit of models
anova(m0, m1, m2, m3)

## Data: d
## Models:
## m0: cbind(totalNFreq, totalEFreq) ~ 1 + (1 | Set) + (1 | EngLexeme)
## m1: cbind(totalNFreq, totalEFreq) ~ (1 | Set) + (1 | EngLexeme) + NormDistanceChosen
## m2: cbind(totalNFreq, totalEFreq) ~ (1 | Set) + (1 | EngLexeme) + NormDistanceChosen + I(NormDist
## m3: cbind(totalNFreq, totalEFreq) ~ (1 | Set) + (1 | EngLexeme) + NormDistanceChosen + I(NormDist
##      npar   AIC    BIC logLik deviance  Chisq Df Pr(>Chisq)
## m0      3 34838 34854 -17416    34832
## m1      4 34817 34839 -17405    34809 22.605  1 1.990e-06 ***
## m2      5 34727 34754 -17358    34717 92.365  1 < 2.2e-16 ***
## m3      6 34677 34710 -17332    34665 51.710  1 6.434e-13 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Summary of final model:
summary(m3)

## Generalized linear mixed model fit by maximum likelihood (Laplace
## Approximation) [glmerMod]
## Family: binomial ( logit )
## Formula: cbind(totalNFreq, totalEFreq) ~ (1 | Set) + (1 | EngLexeme) +
##          NormDistanceChosen + I(NormDistanceChosen^2) + I(NormDistanceChosen^3)
## Data: d
##
##           AIC          BIC    logLik deviance df.resid
## 34677.1    34709.8 -17332.5   34665.1      1737
##
## Scaled residuals:
##      Min       1Q   Median       3Q      Max
## -15.0151  -2.4880  -0.5854   0.8175  24.4700
##
## Random effects:
##  Groups      Name              Variance Std.Dev.
##  EngLexeme (Intercept) 1.918      1.385
##  Set          (Intercept) 1.493      1.222
## Number of obs: 1743, groups: EngLexeme, 135; Set, 67
##
## Fixed effects:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)      0.6460    0.2248   2.873  0.00406 **
## NormDistanceChosen -3.0410    1.0795 -2.817  0.00485 **
```

```
## I(NormDistanceChosen^2) 14.2411      2.7978    5.090 3.58e-07 ***
## I(NormDistanceChosen^3) -13.4971     1.9893   -6.785 1.16e-11 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Correlation of Fixed Effects:
##              (Intr) NrmDsC I(NDC^2
## NrmDstncChs  -0.379
## I(NrmDsC^2)   0.354 -0.983
## I(NrmDsC^3)  -0.334  0.947 -0.988
pHistoricalTotal = plot_model(m3,'eff',
  terms="NormDistanceChosen[all]")
pHistoricalTotal
```



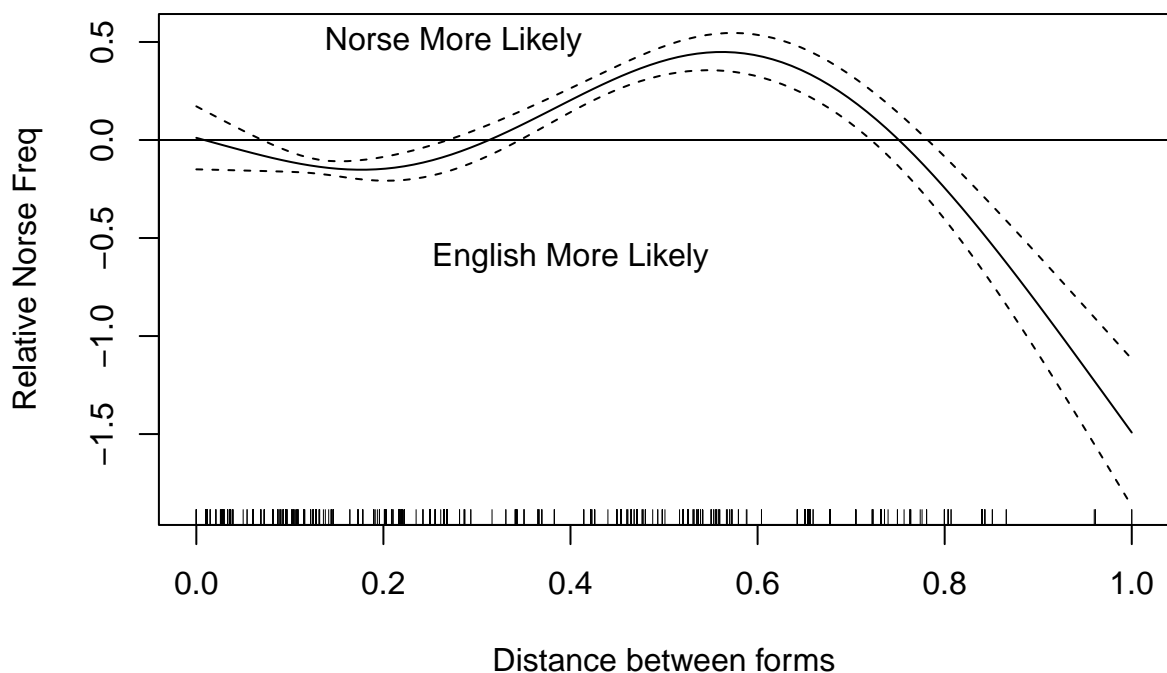
Statistics: (Linear: beta = -3.041, z = -2.817, Wald p = 0.005, LLDiff = 11.3, df = 1, p < 0.001; Quadratic: beta = 14.241, z = 5.09, Wald p < 0.001, LLDiff = 46.2, df = 1, p < 0.001; Cubic: beta = -13.497, z = -6.785, Wald p < 0.001, LLDiff = 25.9, df = 1, p < 0.001)

Same analysis using GAM with binomial distribution:

```
mGAM = gam(cbind(totalNFreq,totalEFreq) ~
  s(NormDistanceChosen,k=4) +
  s(Set,bs="re") +
  s(EngLexeme,bs="re"),
  data = d, family = "binomial")
summary(mGAM)
```

```
##
## Family: binomial
## Link function: logit
##
## Formula:
## cbind(totalNFreq, totalEFreq) ~ s(NormDistanceChosen, k = 4) +
```

```
##      s(Set, bs = "re") + s(EngLexeme, bs = "re")
##
## Parametric coefficients:
##           Estimate Std. Error z value Pr(>|z|)
## (Intercept)  0.5765      0.2331   2.473  0.0134 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Approximate significance of smooth terms:
##           edf Ref.df   Chi.sq p-value
## s(NormDistanceChosen)  2.999     3    138.2 <2e-16 ***
## s(Set)                  44.718    66 741073.2  0.5201
## s(EngLexeme)            81.985   134 689609.4  0.0628 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## R-sq.(adj) =  0.411   Deviance explained = 43.3%
## UBRE = 16.554   Scale est. = 1           n = 1743
plot.gam(mGAM,
  ylab="Relative Norse Freq",
  xlab="Distance between forms",
  select = 1)
abline(h=0)
text(0.275,0.5,"Norse More Likely")
text(0.4,-0.6,"English More Likely")
```

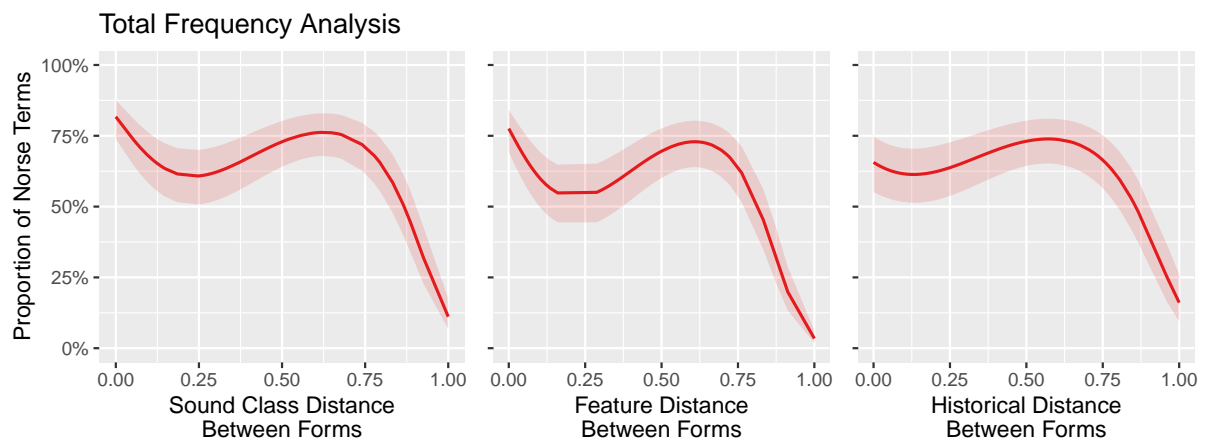


Summary of results for total frequency analyses:

```
bigPlot = grid.arrange(
  pSimpleTotal +
    ggtitle("Total Frequency Analysis") +
    coord_cartesian(ylim = c(0,1))+
    xlab("Sound Class Distance\nBetween Forms")+
    ylab("Proportion of Norse Terms"),
  pFeatureTotal+
    ggtitle("") +
    theme(axis.title.y = element_blank(),
          axis.text.y=element_blank())+
    coord_cartesian(ylim = c(0,1))+
    xlab("Feature Distance\nBetween Forms"),
  pHistoricalTotal+
    ggtitle("") +
    theme(axis.title.y = element_blank(),
          axis.text.y=element_blank())+
    coord_cartesian(ylim = c(0,1))+
    xlab("Historical Distance\nBetween Forms"),
  nrow=1,widths=c(1.3,1,1))

pdf("../results/BigEffectsPlot_totalFreq.pdf",width = 8,height=3)
plot(bigPlot)
dev.off()
```

```
## pdf
## 2
```



5.4 Exploratory analyses

5.4.1 Word Class

It is possible that lexical choices differ by word class. For example, verbs may be more resistant to integration. We test this by adding word class as a predictor. Since some word forms appear as several classes, we treat word class as a series of independent binary variables indicating the possibility or not of the word appearing as this class.

```
classes = unique(trimws(unlist(strsplit(unique(d$Class),"/"))))
classes = classes[!is.na(classes)]
classNames = paste0("Class.",gsub(" ",".",classes))
for(i in 1:length(classes)){
  d[,classNames[i]] = grepl(paste0("^",classes[i]),d$Class)
}
```

Build a null model with non-linear distance, as above:

```
d$NormDistanceChosen = d$NormDistance
mC0 = glmer(cbind(totalNFreq,totalEFreq) ~
  1 + NormDistanceChosen +
  I(NormDistanceChosen^2) +
  I(NormDistanceChosen^3) +
  (1|Set) + (1|EngLexeme),
  data = d, family = "binomial")
```

Adding the verb class variable does not improve the fit of the model:

```
mC1a = update(mC0,~.+Class.verb)
anova(mC0,mC1a)
```

```
## Data: d
## Models:
## mC0: cbind(totalNFreq, totalEFreq) ~ 1 + NormDistanceChosen + I(NormDistanceChosen^2) + I(NormDis
## mC1a: cbind(totalNFreq, totalEFreq) ~ NormDistanceChosen + I(NormDistanceChosen^2) + I(NormDistan
##      npar   AIC    BIC logLik deviance  Chisq Df Pr(>Chisq)
## mC0      6 34557 34589 -17272    34545
## mC1a     7 34558 34597 -17272    34544 0.2408  1    0.6236
```

Adding all other classes does not improve the fit of the model:

```
mC1b = update(mC0,~.+Class.adverb + Class.preposition + Class.noun +
  Class.numeral + Class.adjective + Class.verb +
  Class.indefinite.pronoun + Class.interjection)
```

```
## fixed-effect model matrix is rank deficient so dropping 2 columns / coefficients
## Warning in checkConv(attr("opt", "derivs"), opt$par, ctrl = control$checkConv, :
## Model failed to converge with max|grad| = 0.00857171 (tol = 0.002, component 1)
anova(mC0,mC1b)
```

```
## Data: d
## Models:
## mC0: cbind(totalNFreq, totalEFreq) ~ 1 + NormDistanceChosen + I(NormDistanceChosen^2) + I(NormDis
## mC1b: cbind(totalNFreq, totalEFreq) ~ NormDistanceChosen + I(NormDistanceChosen^2) + I(NormDistan
##      npar   AIC    BIC logLik deviance  Chisq Df Pr(>Chisq)
## mC0      6 34557 34589 -17272    34545
## mC1b    12 34561 34627 -17269    34537 7.1531  6    0.3069
summary(mC1b)
```

```
## Generalized linear mixed model fit by maximum likelihood (Laplace
## Approximation) [glmerMod]
## Family: binomial ( logit )
```



```

## Formula:
## cbind(totalNFreq, totaleFreq) ~ NormDistanceChosen + I(NormDistanceChosen^2) +
##   I(NormDistanceChosen^3) + (1 | Set) + (1 | EngLexeme) + Class.adverb +
##   Class.preposition + Class.noun + Class.numeral + Class.adjective +
##   Class.verb + Class.indefinite.pronoun + Class.interjection
## Data: d
##
##      AIC      BIC   logLik deviance df.resid
## 34561.4 34627.0 -17268.7 34537.4      1731
##
## Scaled residuals:
##      Min       1Q   Median       3Q      Max
## -15.0652  -2.4832  -0.5431   0.7929  24.4495
##
## Random effects:
## Groups      Name      Variance Std.Dev.
## EngLexeme (Intercept) 1.941    1.393
## Set          (Intercept) 1.199    1.095
## Number of obs: 1743, groups: EngLexeme, 135; Set, 67
##
## Fixed effects:
##
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)      1.12406    0.77703   1.447    0.148
## NormDistanceChosen      -9.96860    0.96716 -10.307 <2e-16 ***
## I(NormDistanceChosen^2)      29.84141    2.25590  13.228 <2e-16 ***
## I(NormDistanceChosen^3)     -23.10983    1.54510 -14.957 <2e-16 ***
## Class.adverbTRUE          0.93493    0.90955   1.028    0.304
## Class.nounTRUE            0.20819    0.82897   0.251    0.802
## Class.numeralTRUE         1.84262    1.18738   1.552    0.121
## Class.adjectiveTRUE       -0.46840    0.91586  -0.511    0.609
## Class.verbTRUE            0.05507    0.83885   0.066    0.948
## Class.indefinite.pronounTRUE 1.00817    1.67757   0.601    0.548
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Correlation of Fixed Effects:
##              (Intr) NrmDsC I(NDC^2 I(NDC^3 Clss.dvTRUE Clss.nnTRUE Clss.nmTRUE
## NrmDstncChs -0.131
## I(NrmDsC^2)  0.117 -0.971
## I(NrmDsC^3) -0.104  0.915 -0.982
## Clss.dvTRUE -0.831 -0.011  0.009 -0.007
## Clss.nnTRUE -0.917 -0.002  0.005 -0.006  0.779
## Clss.nmTRUE -0.635  0.011 -0.010  0.008  0.539  0.593
## Clss.djTRUE -0.830 -0.003  0.003 -0.003  0.706  0.783  0.537
## Clss.vrTRUE -0.860 -0.001  0.004 -0.004  0.730  0.804  0.556
## Clss.n.TRUE -0.453 -0.012  0.015 -0.018  0.386  0.425  0.294
##              Clss.djTRUE Clss.vTRUE
## NrmDstncChs
## I(NrmDsC^2)
## I(NrmDsC^3)
## Clss.dvTRUE
## Clss.nnTRUE
## Clss.nmTRUE
## Clss.djTRUE
## Clss.vrTRUE  0.728
## Clss.n.TRUE  0.385      0.398
## fit warnings:
## fixed-effect model matrix is rank deficient so dropping 2 columns / coefficients

```

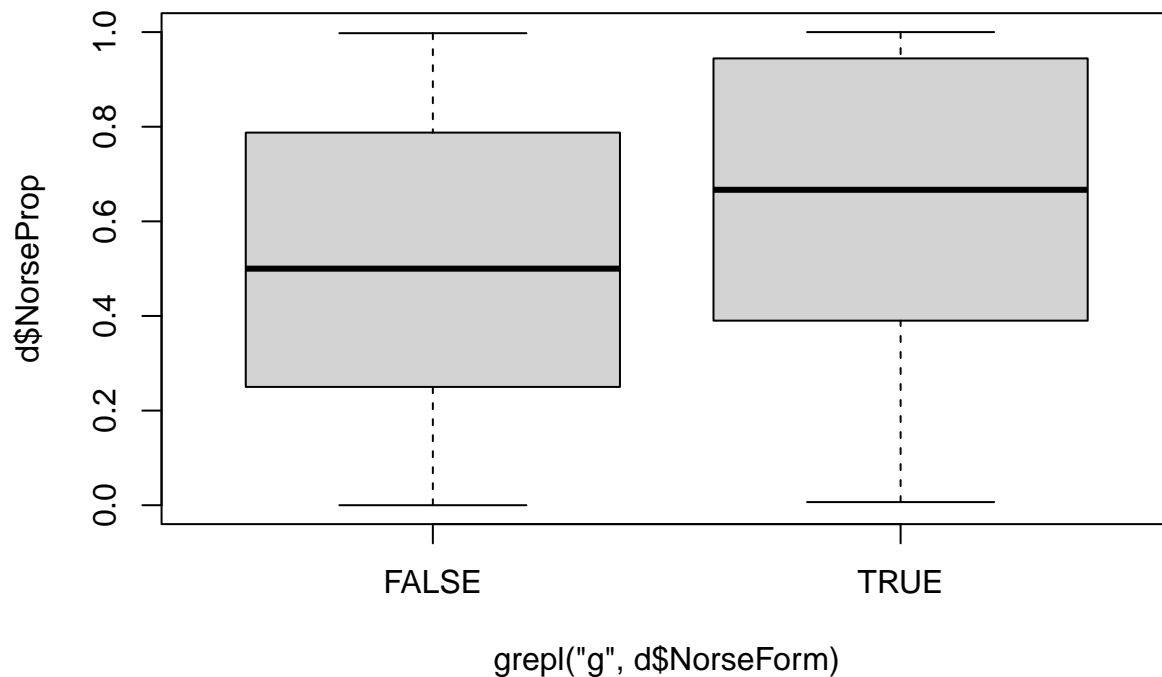
```
## optimizer (Nelder_Mead) convergence code: 0 (OK)
## Model failed to converge with max|grad| = 0.00857171 (tol = 0.002, component 1)
```

5.4.2 Diagnostic consonants

Some Norse words have more salient diagnostic segments - they sound ‘more Norse’. For example, the presence of Velar consonants (e.g. /g/ vs /j/, /sk/ vs post-alveolar fricative, /k/ vs post-alveolar affricate).

For example, the average Norse proportion is higher for Norse forms words with ‘g’:

```
boxplot(d$NorseProp ~ grepl("g",d$NorseForm))
```

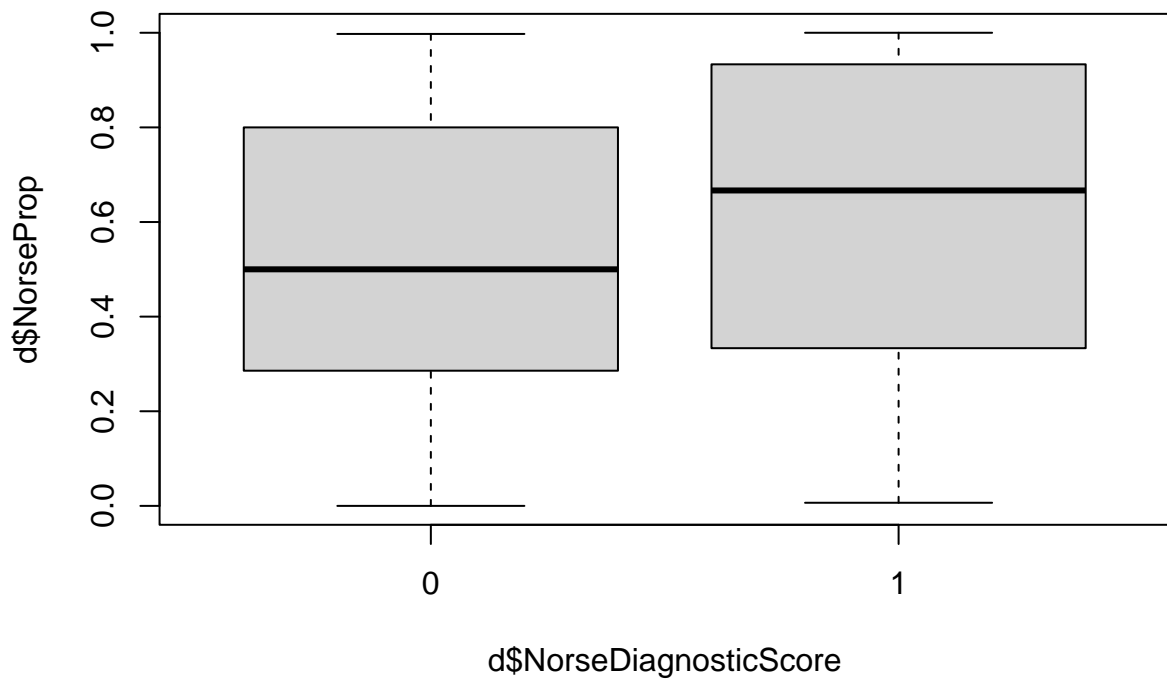


To test this systematically, we calculate a Norse Diagnostic Score: cases where these segments are diagnostic of the differences between Norse and English forms (the score is potentially more than one, but in our data there are only zero and one):

```
d$NorseDiagnosticScore =
  (grepl("g",d$NorseForm) & !grepl("g",d$EngForm))+
  (grepl(" ",d$NorseForm) & !grepl(" ",d$EngForm)) +
  (grepl("sk",d$NorseForm) & !grepl("sk",d$EngForm))+
  # (k in Norse vs in English)
  (grepl("k",d$NorseForm) & grepl(" ",d$EngForm))

d$NorseDiagnosticScore = factor(d$NorseDiagnosticScore)

boxplot(d$NorseProp~ d$NorseDiagnosticScore)
```



```
t.test(d$NorseProp~ d$NorseDiagnosticScore)
```

```
##
##  Welch Two Sample t-test
##
## data:  d$NorseProp by d$NorseDiagnosticScore
## t = -4.2466, df = 1618.4, p-value = 2.294e-05
## alternative hypothesis: true difference in means between group 0 and group 1 is not equal to 0
## 95 percent confidence interval:
##  -0.09465675 -0.03484287
## sample estimates:
## mean in group 0 mean in group 1
##      0.5430542      0.6078040
```

Adding the Norse Diagnostic Score to the model significantly improves the fit of the model. The model suggests that Norse forms which are more obviously Norse have a lower chance of being selected overall.

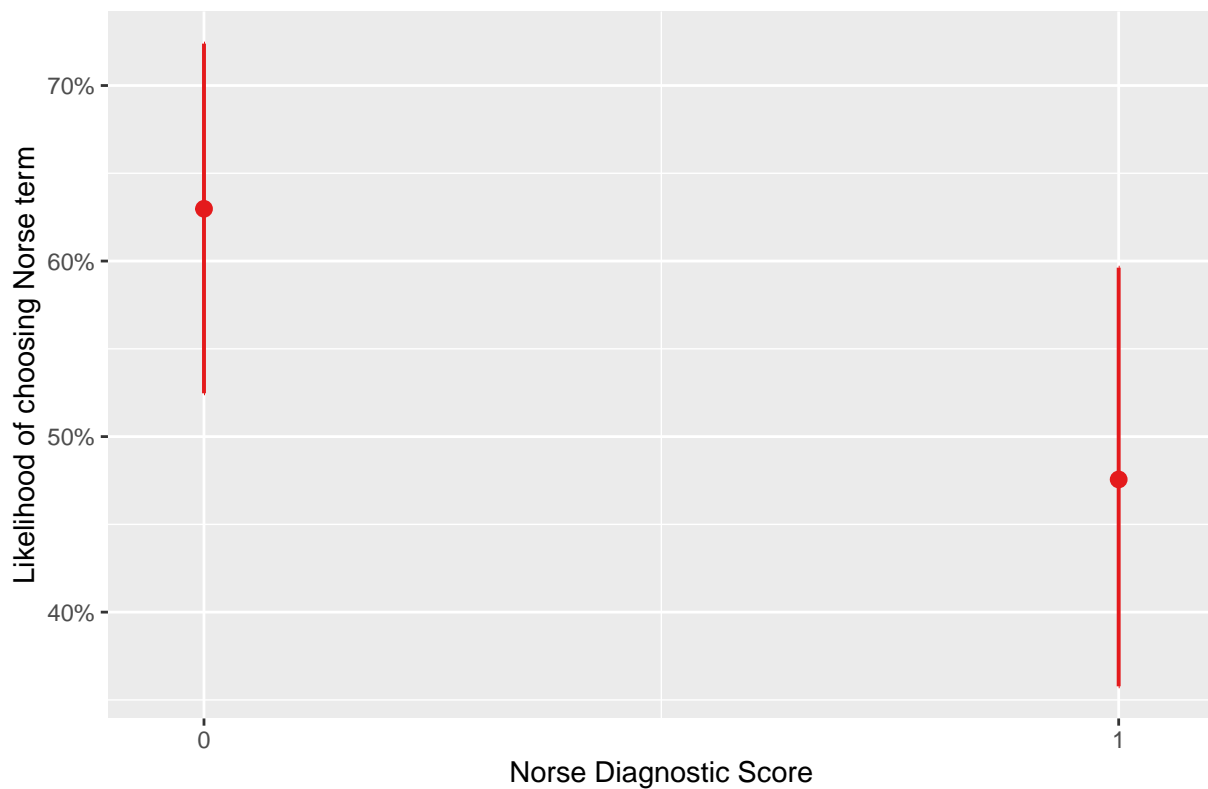
```
d$NormDistanceChosen = d$NormDistance
mD0 = glmer(cbind(totalNFreq,totalEFreq) ~
            1 + NormDistanceChosen +
              I(NormDistanceChosen^2) +
              I(NormDistanceChosen^3) +
              (1|Set) + (1|EngLexeme),
            data = d, family = "binomial")
mD1 = update(mD0, ~.+NorseDiagnosticScore)
anova(mD0,mD1)
```

```
## Data: d
## Models:
## mD0: cbind(totalNFreq, totalEFreq) ~ 1 + NormDistanceChosen + I(NormDistanceChosen^2) + I(NormDis
## mD1: cbind(totalNFreq, totalEFreq) ~ NormDistanceChosen + I(NormDistanceChosen^2) + I(NormDistanc
##      npar   AIC   BIC logLik deviance  Chisq Df Pr(>Chisq)
## mD0      6 34557 34589 -17272   34545
## mD1      7 34542 34580 -17264   34528 16.921  1 3.897e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
summary(mD1)
```

```
## Generalized linear mixed model fit by maximum likelihood (Laplace
## Approximation) [glmerMod]
## Family: binomial ( logit )
## Formula:
## cbind(totalNFreq, totaleFreq) ~ NormDistanceChosen + I(NormDistanceChosen^2) +
## I(NormDistanceChosen^3) + (1 | Set) + (1 | EngLexeme) + NorseDiagnosticScore
## Data: d
##
##      AIC      BIC   logLik deviance df.resid
## 34541.7 34579.9 -17263.8 34527.7    1736
##
## Scaled residuals:
##      Min       1Q   Median       3Q      Max
## -15.0594  -2.4802  -0.5357   0.8251  24.4509
##
## Random effects:
## Groups      Name      Variance Std.Dev.
## EngLexeme (Intercept) 1.963    1.401
## Set      (Intercept) 1.741    1.320
## Number of obs: 1743, groups: EngLexeme, 135; Set, 67
##
## Fixed effects:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)      1.5168    0.2436   6.225 4.81e-10 ***
## NormDistanceChosen -9.9064    0.9445 -10.488 < 2e-16 ***
## I(NormDistanceChosen^2) 29.6575    2.2014 13.472 < 2e-16 ***
## I(NormDistanceChosen^3) -22.9241    1.5105 -15.177 < 2e-16 ***
## NorseDiagnosticScore1 -0.6290    0.1547  -4.066 4.79e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Correlation of Fixed Effects:
##              (Intr) NrmDsC I(NDC^2 I(NDC^3
## NrmDstncChs -0.417
## I(NrmDsC^2)  0.376 -0.969
## I(NrmDsC^3) -0.337  0.910 -0.981
## NrsDgnstcS1 -0.138 -0.011  0.015  -0.020

plot_model(mD1,'eff', terms="NorseDiagnosticScore") +
  ylab("Likelihood of choosing Norse term") +
  xlab("Norse Diagnostic Score") +
  ggtitle("")
```



```
get_model_data(mD1, 'eff', terms="NorseDiagnosticScore")
```

```
## # Predicted probabilities of cbind(totalNFreq, totalEFreq)
```

```
##
```

```
## NorseDiagnosticScore | Predicted | 95% CI | group_col
```

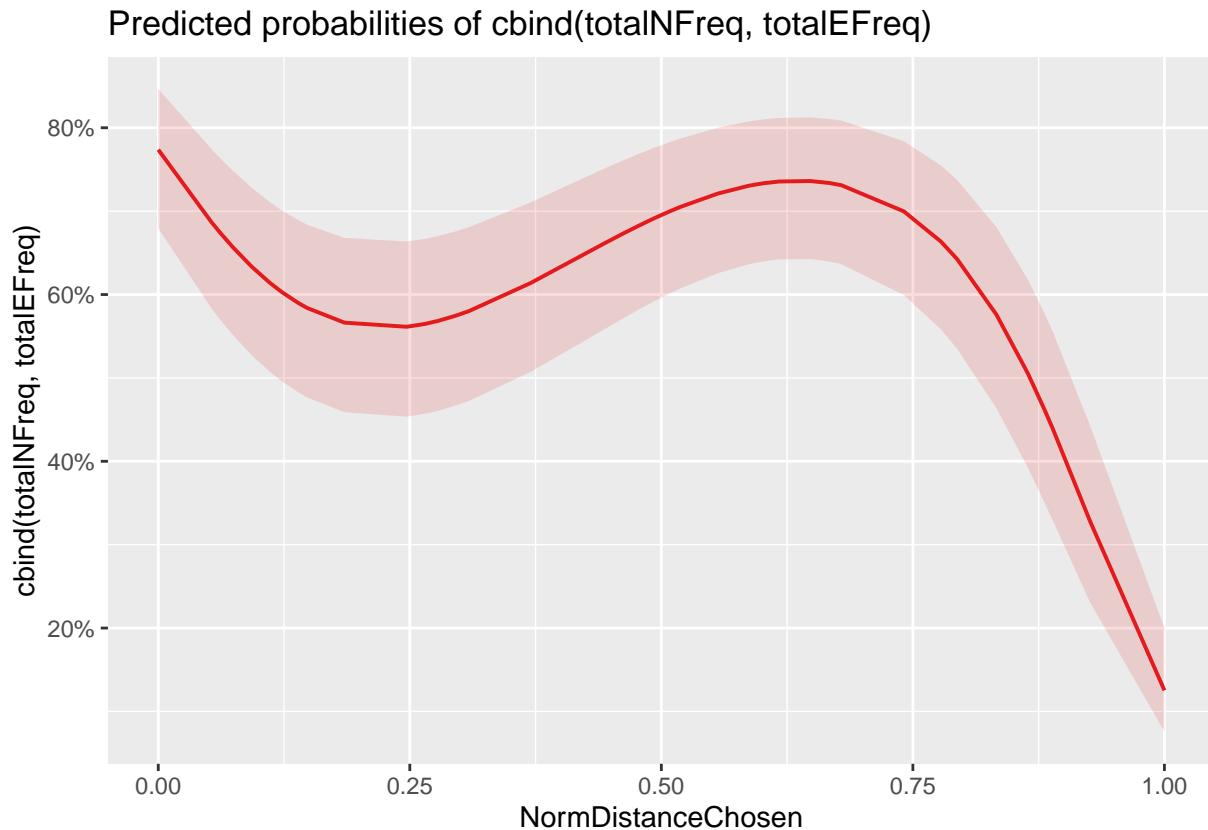
```
## -----
```

```
##           0 | 0.63 | 0.52, 0.72 | 1
```

```
##           1 | 0.48 | 0.36, 0.60 | 1
```

The effect of form distance is unaffected:

```
plot_model(mD1, 'eff', terms="NormDistanceChosen[all]")
```



6 Source-level frequency analysis

This analysis uses source-level observations: each observation is the frequency of a pair of forms within a particular source. This also lets us control for source-level features like the source itself as a random effect, the age of the source, and whether alliteration might affect decisions.

There may be multiple orthographic forms for each unique full form for comparison, so first, we collapse the data over unique pairs. We also restructure the data to be in 'long' format.

```
d2 = data.frame()
for(i in 1:length(englishFrequencyColumns)){
  dx = data.frame(
    Set = d$Set,
    EngForm = d$EngForm,
    NorseForm = d$NorseForm,
    NormDistance = d$NormDistance,
    NormFeatureDistance = d$NormFeatureDistance,
    NormHistoricalDistance = d$NormHistoricalDistance,
    NFreq = d[,norseFrequencyColumns[i]],
    EFreq = d[,englishFrequencyColumns[i]],
    Source = englishFrequencyColumns[i],
    Alliteration = d$Alliteration)
  dx = dx[rowSums(dx[,c("NFreq", "EFreq")], na.rm = T) > 0,]
  uforms = paste(dx$Set, dx$EngForm, dx$NorseForm)
  dx = data.frame(
    Set = tapply(dx$Set, uforms, head, n=1),
    EngForm = tapply(dx$EngForm, uforms, head, n=1),
    NorseForm = tapply(dx$NorseForm, uforms, head, n=1),
    NormDistance = tapply(dx$NormDistance, uforms, head, n=1),
    NormFeatureDistance = tapply(dx$NormFeatureDistance, uforms, head, n=1),
    NormHistoricalDistance = tapply(dx$NormHistoricalDistance, uforms, head, n=1),
    NFreq = tapply(dx$NFreq, uforms, sum, na.rm=T),
```

```

  EFreq = tapply(dx$EFreq,uforms,sum,na.rm=T),
  Source = englishFrequencyColumns[i],
  Alliteration = tapply(dx$Alliteration,uforms,head,n=1))
d2 = rbind(d2,dx)
}

d2$Set = factor(d2$Set)
d2$Source = factor(d2$Source)

# Make a variable for each unique form within a set
d2$EngForm2 = paste(d2$Set,d2$EngForm)
d2$NorseForm2 = paste(d2$Set,d2$NorseForm)
d2$EngForm2 = factor(d2$EngForm2)

# Remove any cases with zero frequency
d2 = d2[!(d2$NFreq==0 & d2$EFreq==0),]

# Proportion of norse forms
d2$NProp = 100*(d2$NFreq/(d2$NFreq+d2$EFreq))

```

Modelling age is difficult, since exact dates are not known. It would be possible to model the date of publication as an ordinal variable, for example assuming FCPC > Ormulum > {Havelok, GenAndEx} > CursorMundi > {Mannyng, GawainPoet} > WarsAlexander > StErkenwald. However, the texts from the Corpus of Middle English span several centuries, and so are hard to place in this order. Instead, we simply use the century as an ordered category.

```

AgeCategories =
  c("EFreqFCPC"=12,
    "EFreqOrmulum"=12,
    "EFreqHavelok"=13,
    "EFreqGenAndEx" = 14,
    "EFreqCursorMundi" = 14,
    "EFreqGawainPoet" = 14,
    "EFreqStErkenwald" = 15,
    "EFreqMannyng" = 15,
    "EFreqWarsAlexander" = 15,
    "EFreqLinc"=15,
    "EFreqNott"=15,
    "EFreqNorf"=15,
    "EFreqRolle"=15)

d2$Age = AgeCategories[d2$Source]
d2$Age = factor(d2$Age,ordered=T)

contrasts(d2$Age) = contr.sum(length(unique(d2$Age)))

mean(d2$NFreq/(d2$NFreq+d2$EFreq))

## [1] 0.5613844
range(d2$NFreq/(d2$NFreq+d2$EFreq))

## [1] 0 1
mean(d2$NormDistance)

## [1] 0.3564722

```

Alliteration only applies to poetry sources: Gawain, St Erkenwold, and Wars of Alexander. So turn all others to 'false':

```

d2[!d2$Source %in%
  c("EFreqGawainPoet",
    "EFreqStErkenwald",
    "EFreqWarsAlexander",
    "EFreqLinc",
    "EFreqNott",
    "EFreqNorf",
    "EFreqRolle"),]$Alliteration = FALSE
d2$Alliteration = factor(d2$Alliteration)

gPairs = ggpairs(d2[,c("NormDistance", "NormFeatureDistance",
  "NormHistoricalDistance", "NProp")],
  columnLabels = c("Sound Class Distance",
    "Feature Distance",
    "Historical Distance",
    "Freq (% Norse)"),
  rowLabels = c("Sound Class Distance",
    "Feature Distance",
    "Historical Distance",
    "Freq (% Norse)"))

pdf("../results/GPairs.pdf", width=6, height=5.5)
gPairs
dev.off()

## pdf
## 2

```

6.1 Simple distance

Use a mixed effects model to predict the frequency of Norse and English forms by a random effect for cognate Set and Source. We add a main effect of age of source, then introduce the normalised distance measure along with its non-linear terms.

```

m0 = glmer(cbind(NFreq, EFreq) ~ Age + Alliteration +
  (1|Set) + (1|Source) +
  (1|EngForm2),
  data = d2, family = "binomial",
  glmerControl(optimizer = "bobyqa"))
# Add the distance measure
m1 = update(m0, ~. + NormDistance)
# Add quadratic term
m2 = update(m1, ~. + I(NormDistance^2))
# Add cubic term
m3 = update(m2, ~. + I(NormDistance^3))

```

Compare fit of models:

```
anova(m0, m1, m2, m3)
```

```

## Data: d2
## Models:
## m0: cbind(NFreq, EFreq) ~ Age + Alliteration + (1 | Set) + (1 | Source) + (1 | EngForm2)
## m1: cbind(NFreq, EFreq) ~ Age + Alliteration + (1 | Set) + (1 | Source) + (1 | EngForm2) + NormDi
## m2: cbind(NFreq, EFreq) ~ Age + Alliteration + (1 | Set) + (1 | Source) + (1 | EngForm2) + NormDi
## m3: cbind(NFreq, EFreq) ~ Age + Alliteration + (1 | Set) + (1 | Source) + (1 | EngForm2) + NormDi
##      npar   AIC    BIC logLik deviance  Chisq Df Pr(>Chisq)
## m0      8 21549 21589 -10766    21533
## m1      9 21369 21414 -10675    21351 182.08  1 < 2.2e-16 ***
## m2     10 21267 21318 -10624    21247 103.45  1 < 2.2e-16 ***

```



```
## m3    11 20994 21050 -10486    20972 274.80 1 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Summary of final model:

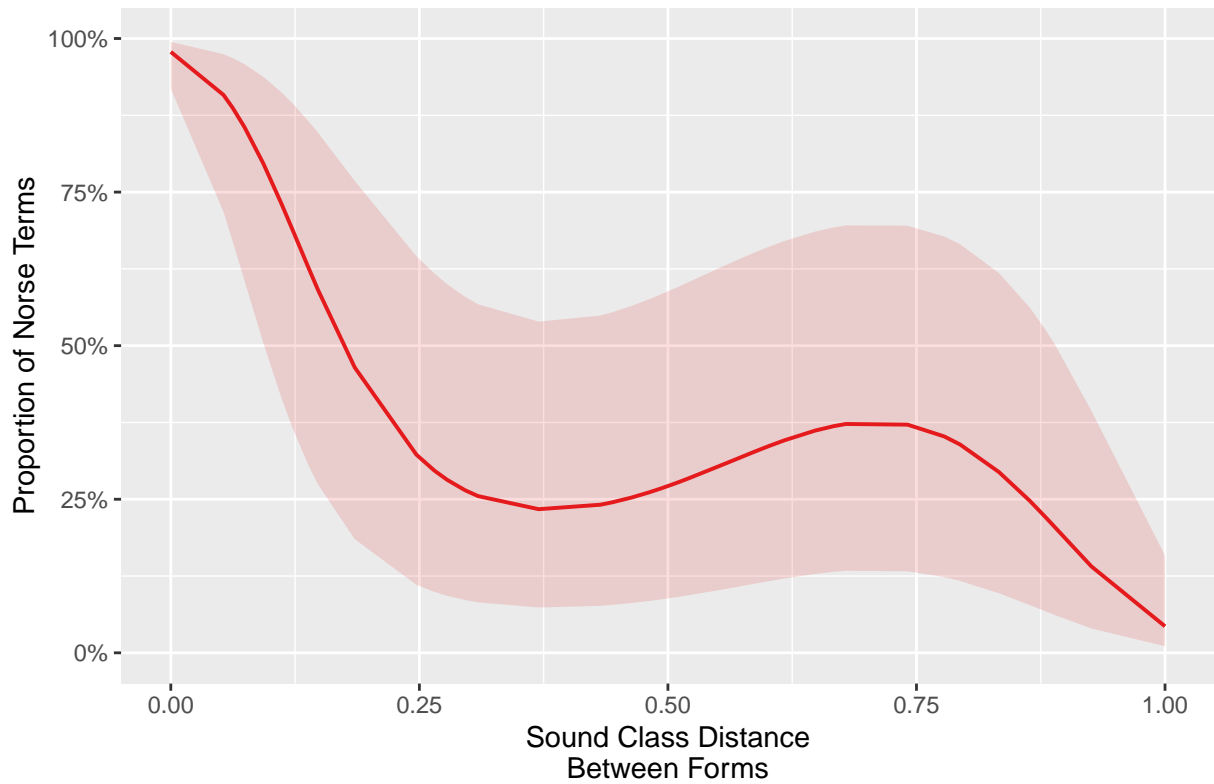
```
summary(m3)
```

```
## Generalized linear mixed model fit by maximum likelihood (Laplace
##   Approximation) [glmerMod]
## Family: binomial ( logit )
## Formula: cbind(NFreq, EFreq) ~ Age + Alliteration + (1 | Set) + (1 | Source) +
##   (1 | EngForm2) + NormDistance + I(NormDistance^2) + I(NormDistance^3)
## Data: d2
## Control: glmerControl(optimizer = "bobyqa")
##
##      AIC      BIC    logLik deviance df.resid
## 20994.5  21050.1 -10486.2  20972.5     1150
##
## Scaled residuals:
##      Min       1Q   Median       3Q      Max
## -98.823  -1.248   0.125   1.133  174.927
##
## Random effects:
## Groups   Name                Variance Std.Dev.
## EngForm2 (Intercept)  4.411      2.100
## Set       (Intercept)  5.306      2.304
## Source   (Intercept)  4.370      2.091
## Number of obs: 1161, groups:  EngForm2, 128; Set, 67; Source, 13
##
## Fixed effects:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    3.48062    0.83091   4.189  2.8e-05 ***
## Age1            1.08270    1.26474   0.856   0.392
## Age2           -1.26811    1.62546  -0.780   0.435
## Age3           -0.01457    1.11874  -0.013   0.990
## AlliterationTRUE  2.28044    0.07102  32.112 < 2e-16 ***
## NormDistance   -31.86647    1.55309 -20.518 < 2e-16 ***
## I(NormDistance^2) 64.16229    3.46394  18.523 < 2e-16 ***
## I(NormDistance^3) -39.19896    2.28369 -17.165 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Correlation of Fixed Effects:
##              (Intr) Age1   Age2   Age3   AlTRUE NrmDst I(ND^2
## Age1              0.012
## Age2              0.397 -0.515
## Age3             -0.183 -0.263 -0.491
## AlltrtnTRUE     -0.012  0.015 -0.010  0.000
## NormDistanc     -0.211  0.001 -0.003  0.002 -0.025
## I(NrmDst^2)      0.191 -0.001  0.003 -0.002  0.026 -0.970
## I(NrmDst^3)     -0.172  0.001 -0.003  0.002 -0.028  0.911 -0.982
```

Norse frequency varies with distance. The plot below shows the marginal effects, holding discrete predictors constant at their proportions (not reference level):

```
pNormSimple =
  plot_model(m3, 'eff', terms="NormDistance[all]") +
  xlab("Sound Class Distance\nBetween Forms") +
  ylab("Proportion of Norse Terms") +
```

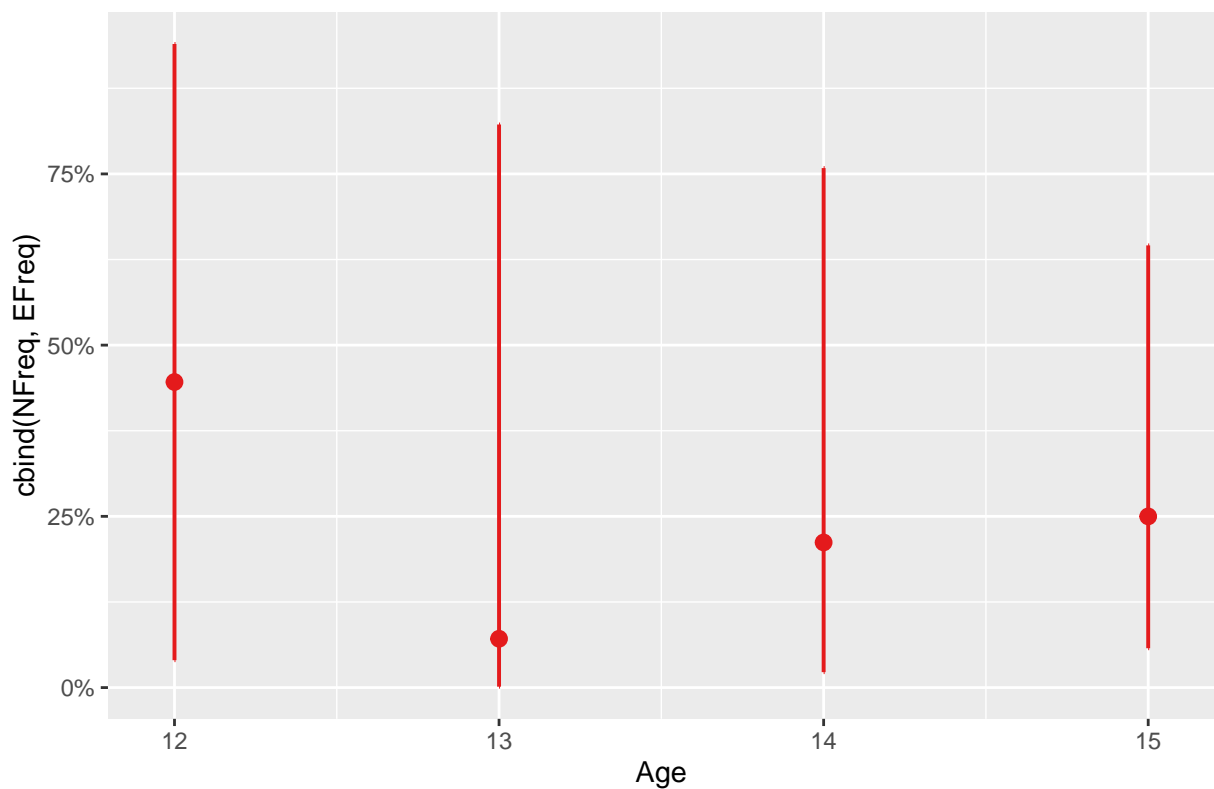
```
ggtitle("")+
  coord_cartesian(ylim=c(0,1))
pNormSimple
```



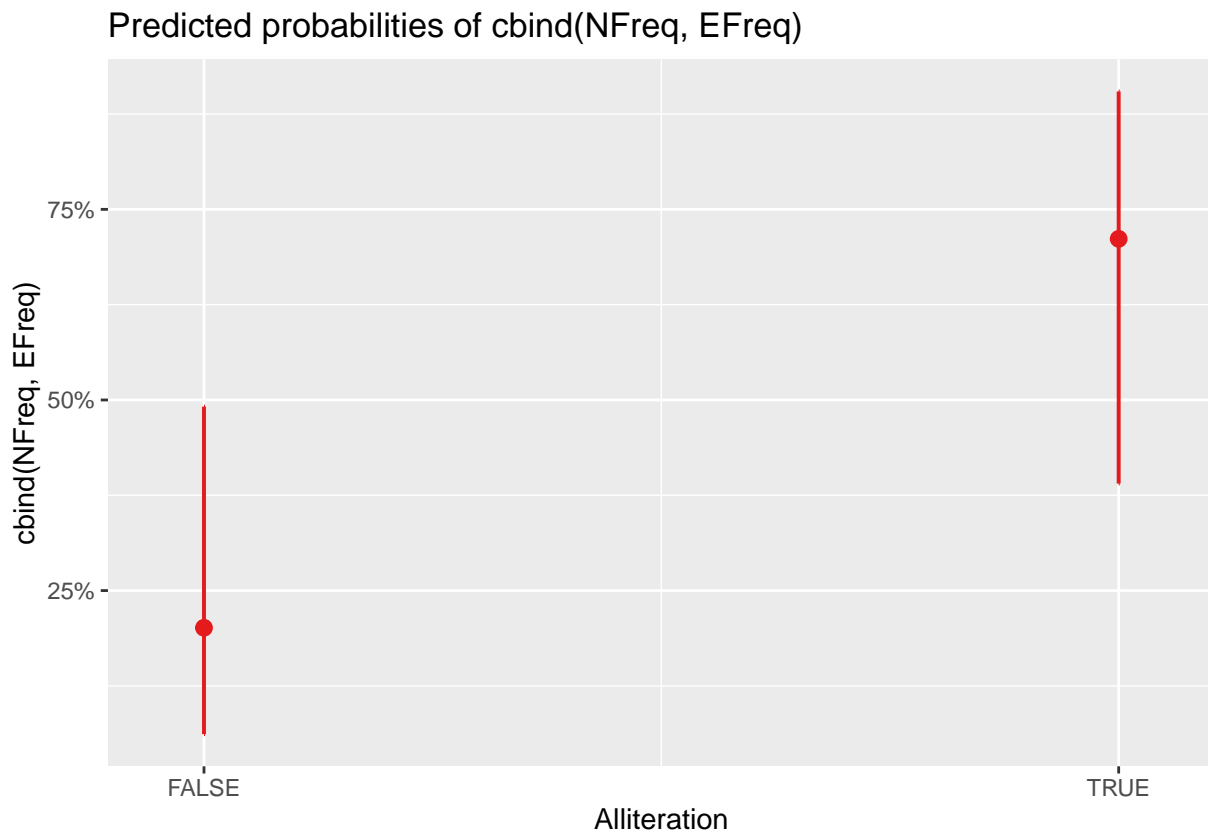
Effects over time and of alliteration:

```
plot_model(m3, 'eff', terms="Age [all]")
```

Predicted probabilities of cbind(NFreq, EFreq)



```
plot_model(m3, 'eff', terms="Alliteration")
```

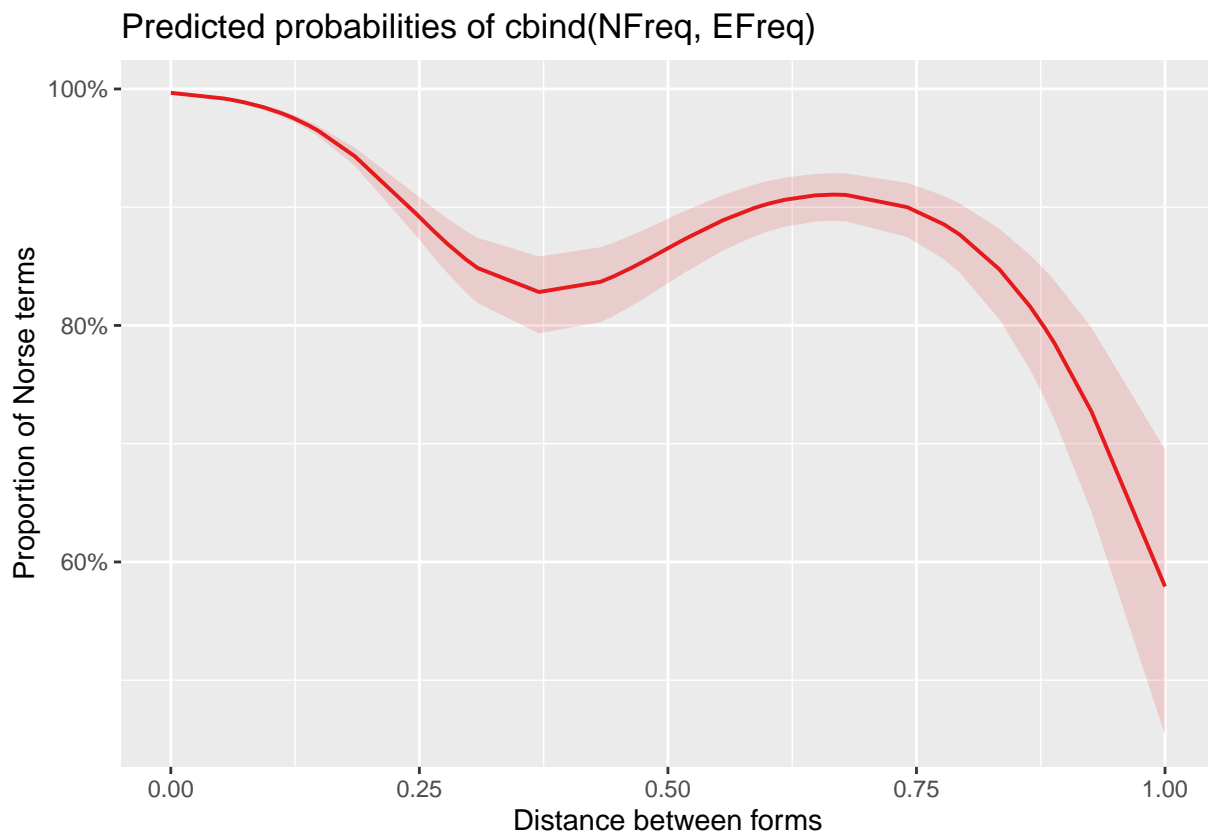


The GAM model seems a little different, showing a smaller effect at lower distances.

```
mGAM = gam(cbind(NFreq, EFreq) ~
            s(NormDistance, k=4) +
            s(Age, bs = "re") +
            s(Set, bs="re") +
            s(Source, bs="re")+
            s(EngForm2, bs="re"),
            data = d2, family = "binomial")
summary(mGAM)
```

```
##
## Family: binomial
## Link function: logit
##
## Formula:
## cbind(NFreq, EFreq) ~ s(NormDistance, k = 4) + s(Age, bs = "re") +
##      s(Set, bs = "re") + s(Source, bs = "re") + s(EngForm2, bs = "re")
##
## Parametric coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -0.04986    2.13567  -0.023   0.981
##
## Approximate significance of smooth terms:
##              edf Ref.df   Chi.sq p-value
## s(NormDistance)  2.99960     3 4.546e+02 < 2e-16 ***
## s(Age)           0.01179     3 2.226e+02 0.99832
## s(Set)          49.17843    66 3.422e+07 0.10349
## s(Source)       11.98614    12 1.177e+07 < 2e-16 ***
## s(EngForm2)     71.28579   127 4.201e+07 0.00559 **
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## R-sq.(adj) =  0.71   Deviance explained = 69.2%
## UBRE = 17.133   Scale est. = 1          n = 1161
plot_model(mGAM,"pred",terms="NormDistance") +
  ylab("Proportion of Norse terms") +
  xlab("Distance between forms")
```



Statistics for the simple measure:

(Linear: $\beta = -31.866$, $z = -20.518$, Wald $p < 0.001$, LLDiff = 91, $df = 1$, $p < 0.001$; Quadratic: $\beta = 64.162$, $z = 18.523$, Wald $p < 0.001$, LLDiff = 51.7, $df = 1$, $p < 0.001$; Cubic: $\beta = -39.199$, $z = -17.165$, Wald $p < 0.001$, LLDiff = 137.4, $df = 1$, $p < 0.001$)

6.2 Feature-based distance

Same analysis as the simple distance above:

```
m0 = glmer(cbind(NFreq,EFreq) ~ Age + Alliteration
           + (1|Set) + (1|Source) +
           (1|EngForm2),
           data = d2, family = "binomial")
m1 = update(m0,~.+NormFeatureDistance)
m2 = update(m1,~.+I(NormFeatureDistance^2))
m3 = update(m2,~.+I(NormFeatureDistance^3))
```

```
## Warning in checkConv(attr("derivs"), opt$par, ctrl = control$checkConv, :
## Model failed to converge with max|grad| = 0.030578 (tol = 0.002, component 1)
```

```
anova(m0,m1,m2,m3)
```

```
## Data: d2
```

```
## Models:
```

```
## m0: cbind(NFreq, EFreq) ~ Age + Alliteration + (1 | Set) + (1 | Source) + (1 | EngForm2)
```

```
## m1: cbind(NFreq, EFreq) ~ Age + Alliteration + (1 | Set) + (1 | Source) + (1 | EngForm2) + NormFe
```

```
## m2: cbind(NFreq, EFreq) ~ Age + Alliteration + (1 | Set) + (1 | Source) + (1 | EngForm2) + NormFe
```

```
## m3: cbind(NFreq, EFreq) ~ Age + Alliteration + (1 | Set) + (1 | Source) + (1 | EngForm2) + NormFe
```

```
##      npar   AIC    BIC logLik deviance  Chisq Df Pr(>Chisq)
```

```
## m0      8 21549 21589 -10766    21533
```

```
## m1      9 21334 21380 -10658    21316 216.84  1 < 2.2e-16 ***
```

```
## m2     10 21226 21276 -10603    21206 110.34  1 < 2.2e-16 ***
```

```
## m3     11 21049 21105 -10514    21027 178.35  1 < 2.2e-16 ***
```

```
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
summary(m3)
```

```
## Generalized linear mixed model fit by maximum likelihood (Laplace
```

```
## Approximation) [glmerMod]
```

```
## Family: binomial ( logit )
```

```
## Formula: cbind(NFreq, EFreq) ~ Age + Alliteration + (1 | Set) + (1 | Source) +
```

```
## (1 | EngForm2) + NormFeatureDistance + I(NormFeatureDistance^2) +
```

```
## I(NormFeatureDistance^3)
```

```
## Data: d2
```

```
##
```

```
##      AIC      BIC    logLik deviance df.resid
```

```
## 21049.3 21104.9 -10513.6 21027.3      1150
```

```
##
```

```
## Scaled residuals:
```

```
##      Min      1Q  Median      3Q      Max
```

```
## -96.751 -1.293   0.119   1.134 176.287
```

```
##
```

```
## Random effects:
```

```
## Groups   Name      Variance Std.Dev.
```

```
## EngForm2 (Intercept) 4.904    2.214
```

```
## Set      (Intercept) 5.691    2.386
```

```
## Source   (Intercept) 4.354    2.087
```

```
## Number of obs: 1161, groups: EngForm2, 128; Set, 67; Source, 13
```

```
##
```

```
## Fixed effects:
```

```
##              Estimate Std. Error z value Pr(>|z|)
```

```
## (Intercept)      1.80918    0.82816   2.185  0.0289 *
```

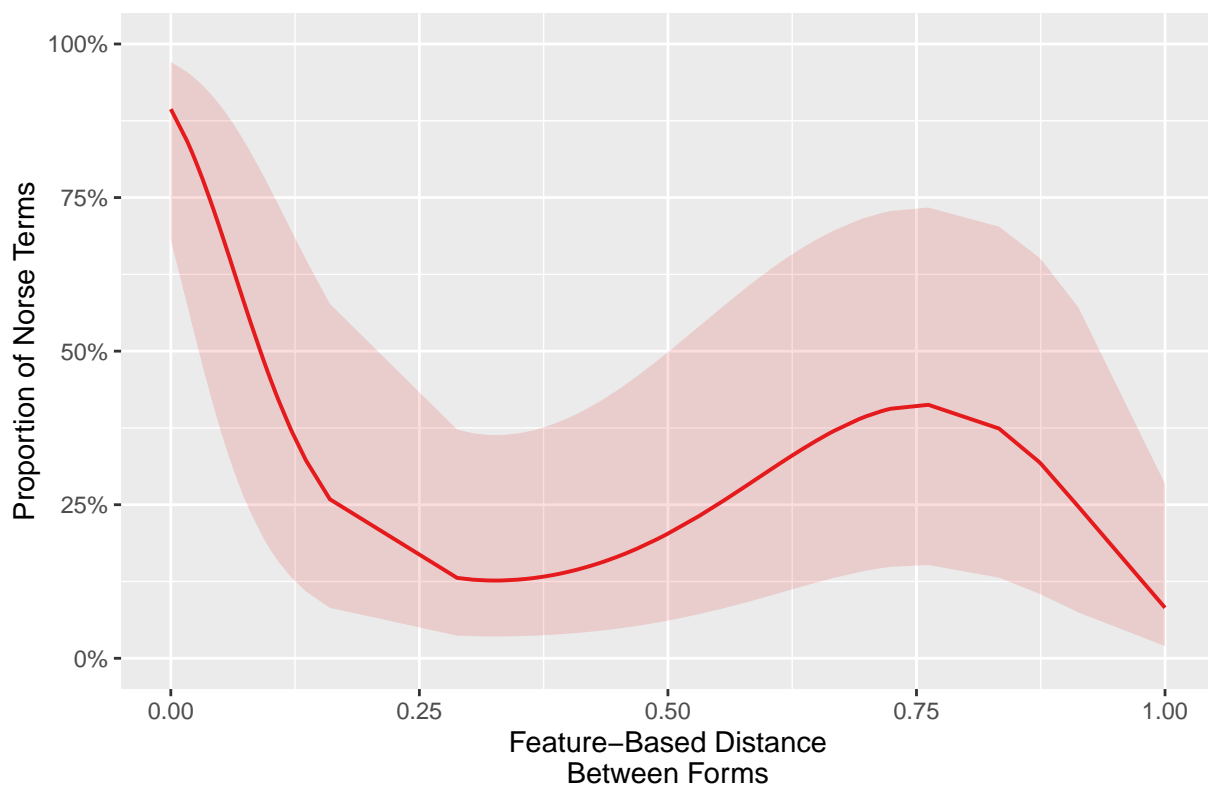
```
## Age1              1.09199    1.27553   0.856  0.3919
```

```
## Age2             -1.26383    1.64662  -0.768  0.4428
```

```
## Age3             -0.01883    1.12408  -0.017  0.9866
```

```
## AlliterationTRUE          2.29756      0.07117  32.285   <2e-16 ***
## NormFeatureDistance      -29.11464     1.79625 -16.209   <2e-16 ***
## I(NormFeatureDistance^2)  63.88467     4.43249  14.413   <2e-16 ***
## I(NormFeatureDistance^3) -39.31342     2.97841 -13.199   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Correlation of Fixed Effects:
##          (Intr) Age1   Age2   Age3   AlTRUE NrmFtD I(NFD^2
## Age1          0.007
## Age2          0.404 -0.520
## Age3          -0.190 -0.259 -0.493
## AlltrtnTRUE  -0.014  0.016 -0.010  0.000
## NrmFtrDstnc  -0.140 -0.002  0.001  0.004 -0.030
## I(NrmFtD^2)   0.131  0.003 -0.001 -0.003  0.033 -0.983
## I(NrmFtD^3)  -0.121 -0.003  0.002  0.003 -0.036  0.940 -0.986
## optimizer (Nelder_Mead) convergence code: 0 (OK)
## Model failed to converge with max|grad| = 0.030578 (tol = 0.002, component 1)
```

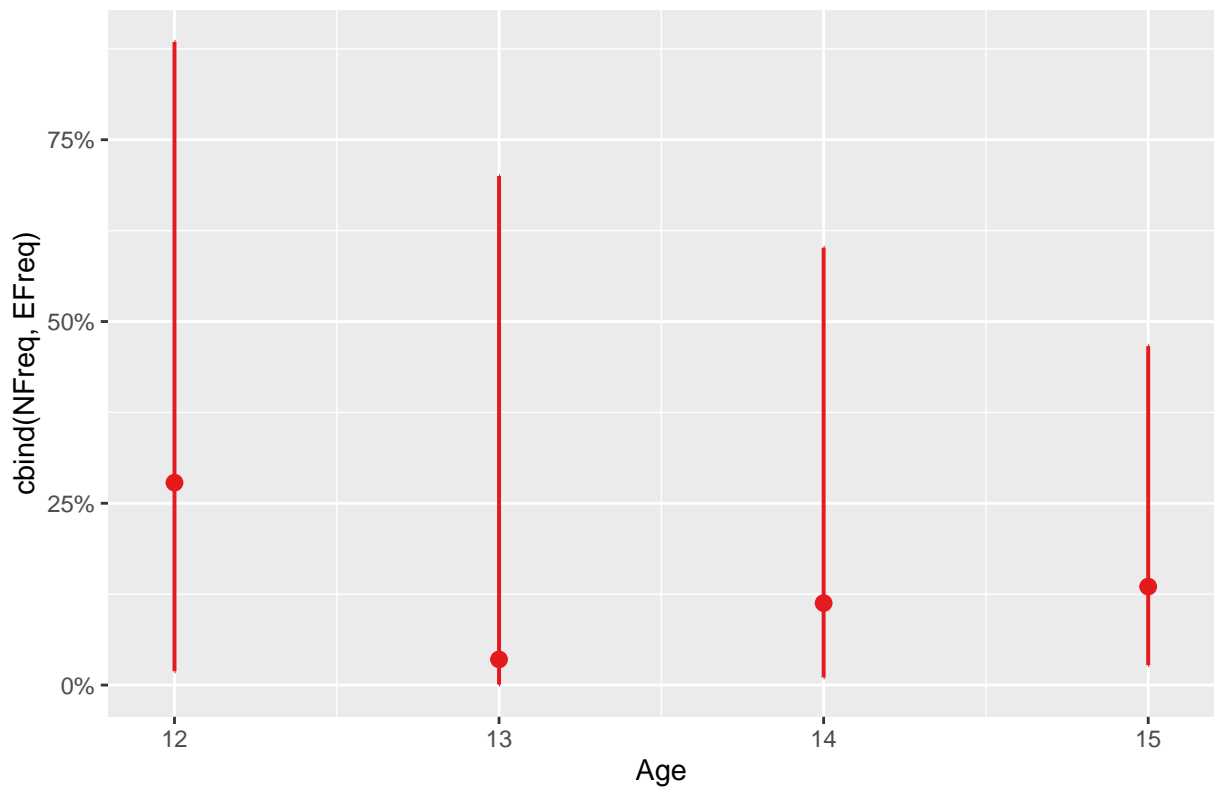
```
pNormFeature =
  plot_model(m3,'eff', terms="NormFeatureDistance[all]")+
    xlab("Feature-Based Distance\nBetween Forms") +
    ylab("Proportion of Norse Terms") +
    ggtitle("")+
    coord_cartesian(ylim=c(0,1))
pNormFeature
```



Effects over time and of alliteration:

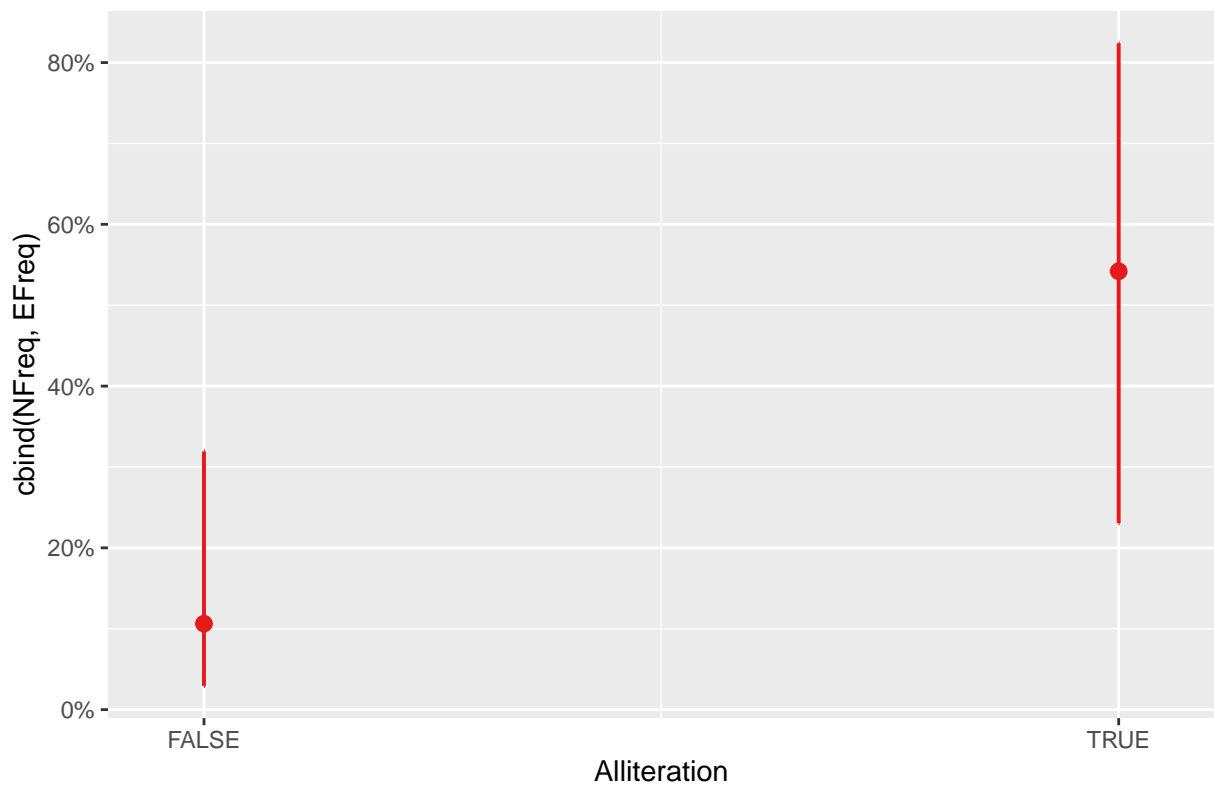
```
plot_model(m3,'eff', terms="Age [all]")
```

Predicted probabilities of cbind(NFreq, EFreq)



```
plot_model(m3, 'eff', terms="Alliteration")
```

Predicted probabilities of cbind(NFreq, EFreq)



GAM model:

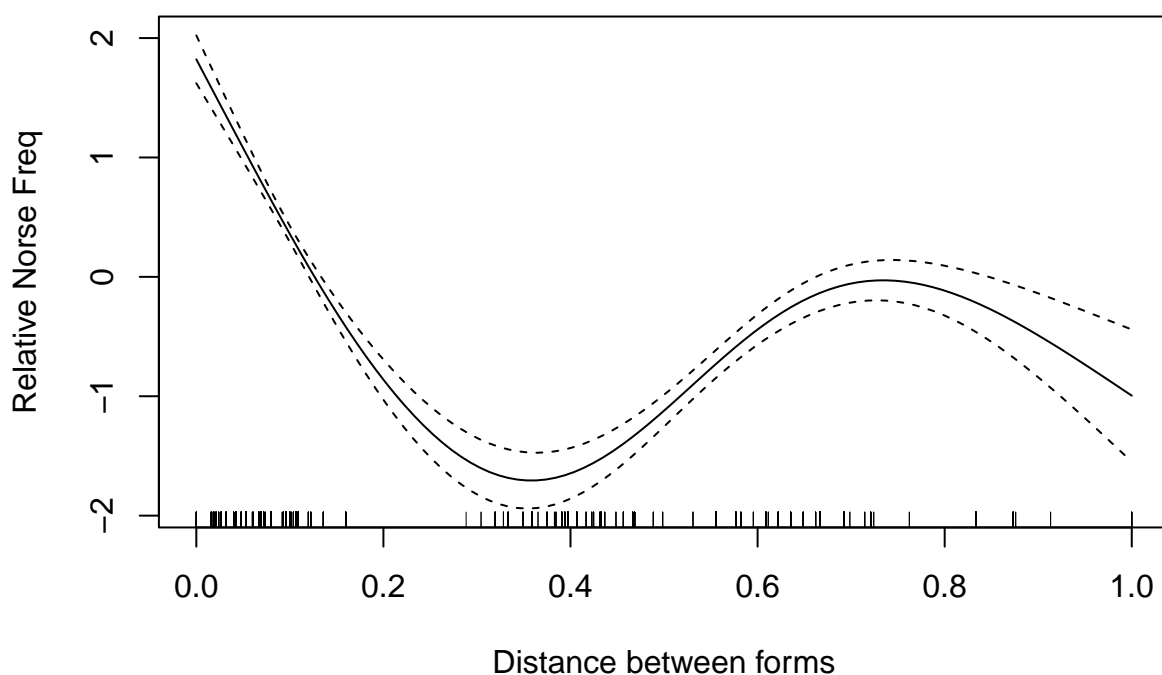
```
mGAM = gam(cbind(NFreq, EFreq) ~
            s(NormFeatureDistance, k=4) +
```

```

s(Age,bs = "re") +
s(Set,bs="re") +
s(Source,bs="re")+
s(EngForm2,bs="re"),
data = d2, family = "binomial")
summary(mGAM)

##
## Family: binomial
## Link function: logit
##
## Formula:
## cbind(NFreq, EFreq) ~ s(NormFeatureDistance, k = 4) + s(Age,
##   bs = "re") + s(Set, bs = "re") + s(Source, bs = "re") + s(EngForm2,
##   bs = "re")
##
## Parametric coefficients:
##               Estimate Std. Error z value Pr(>|z|)
## (Intercept)  0.03653    1.85085   0.02    0.984
##
## Approximate significance of smooth terms:
##               edf Ref.df   Chi.sq p-value
## s(NormFeatureDistance)  2.999739     3 3.636e+02 <2e-16 ***
## s(Age)                  0.001221     3 1.808e+00  0.9979
## s(Set)                  52.222117    66 4.937e+07  0.0659 .
## s(Source)               11.995957    12 1.625e+07 <2e-16 ***
## s(EngForm2)             69.223418   127 1.669e+07  0.2271
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## R-sq.(adj) =  0.711   Deviance explained = 69.1%
## UBRE = 17.188   Scale est. = 1          n = 1161
plot.gam(mGAM,
  ylab="Relative Norse Freq",
  xlab="Distance between forms",
  select = 1)

```



Linear: $\beta = -29.115$, $z = -16.209$, Wald $p < 0.001$, LLDiff = 108.4, $df = 1$, $p < 0.001$; Quadratic: $\beta = 63.885$, $z = 14.413$, Wald $p < 0.001$, LLDiff = 55.2, $df = 1$, $p < 0.001$; Cubic: $\beta = -39.313$, $z = -13.199$, Wald $p < 0.001$, LLDiff = 89.2, $df = 1$, $p < 0.001$

6.3 Historical distance

Same analysis as the simple distance above:

```
m0 = glmer(cbind(NFreq,EFreq) ~ Age + Alliteration +
           (1|Set) + (1|Source) +
           (1|EngForm2),
           data = d2, family = "binomial",
           control=glmerControl(optimizer="bobyqa"))
m1 = update(m0,~.+NormHistoricalDistance)
m2 = update(m1,~.+I(NormHistoricalDistance^2),
           control=glmerControl(optimizer="bobyqa"))
m3 = update(m2,~.+I(NormHistoricalDistance^3),
           control=glmerControl(optimizer="bobyqa"))
anova(m0,m1,m2,m3)

## Data: d2
## Models:
## m0: cbind(NFreq, EFreq) ~ Age + Alliteration + (1 | Set) + (1 | Source) + (1 | EngForm2)
## m1: cbind(NFreq, EFreq) ~ Age + Alliteration + (1 | Set) + (1 | Source) + (1 | EngForm2) + NormHi
## m2: cbind(NFreq, EFreq) ~ Age + Alliteration + (1 | Set) + (1 | Source) + (1 | EngForm2) + NormHi
## m3: cbind(NFreq, EFreq) ~ Age + Alliteration + (1 | Set) + (1 | Source) + (1 | EngForm2) + NormHi
##      npar   AIC    BIC logLik deviance   Chisq Df Pr(>Chisq)
## m0      8 21549 21589 -10766    21533
## m1      9 21361 21406 -10671    21343 190.0816  1 < 2.2e-16 ***
## m2     10 21339 21390 -10660    21319  23.5835  1 1.196e-06 ***
## m3     11 21338 21393 -10658    21316   3.6503  1  0.05606 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

For this measure, the cubic term is marginal.

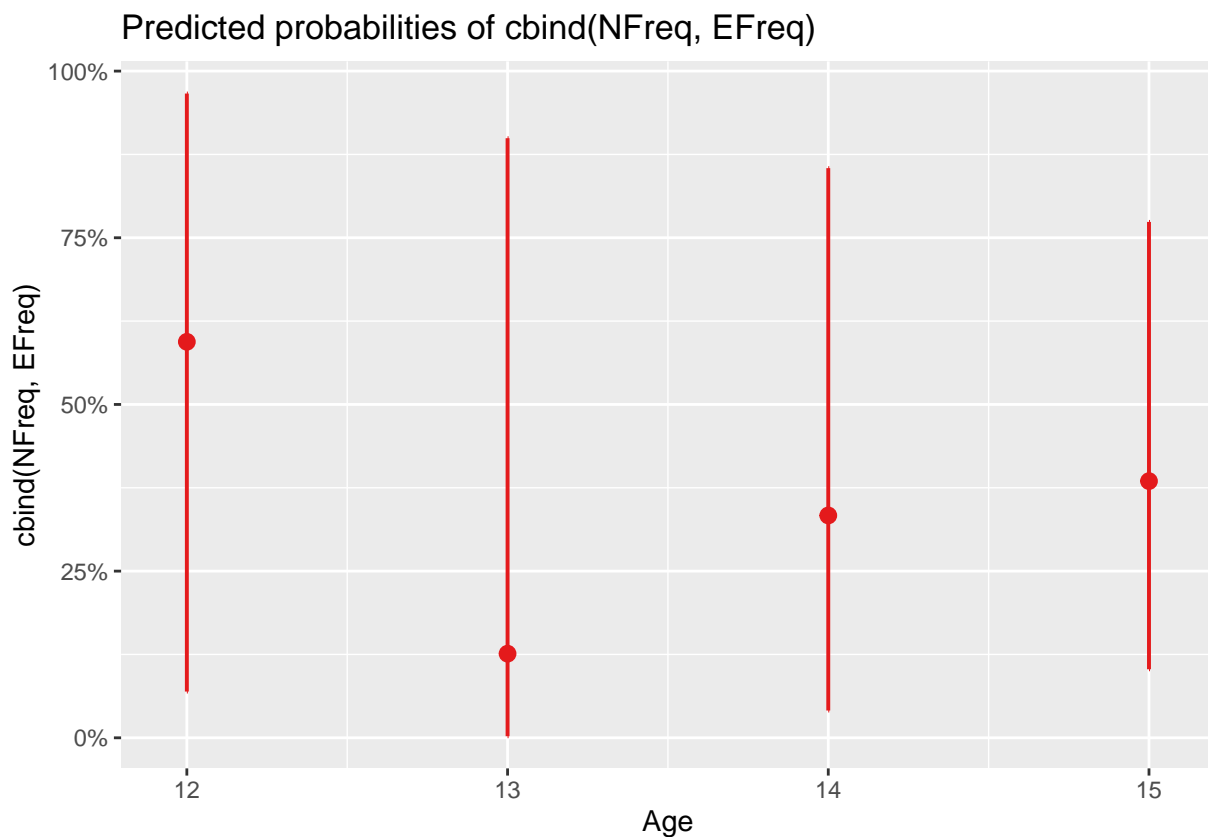
This is the cubic model:

```
summary(m3)

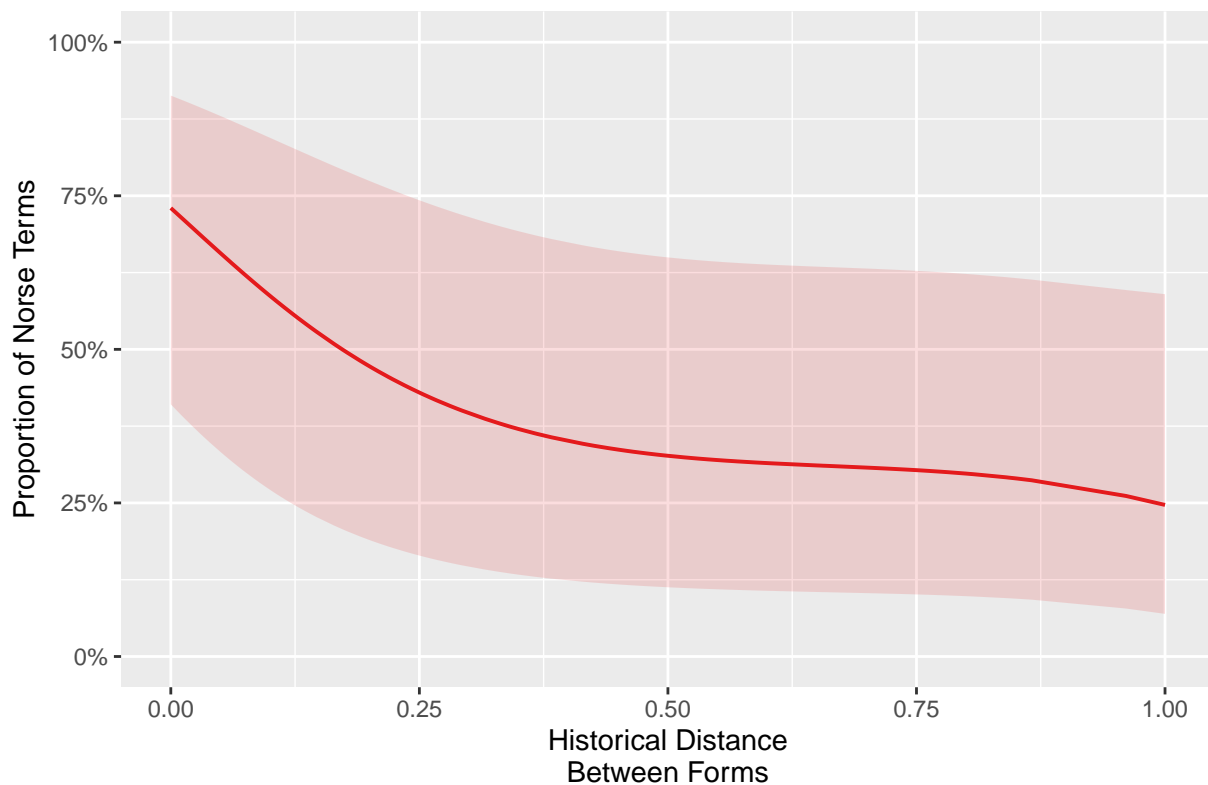
## Generalized linear mixed model fit by maximum likelihood (Laplace
## Approximation) [glmerMod]
## Family: binomial ( logit )
## Formula: cbind(NFreq, EFreq) ~ Age + Alliteration + (1 | Set) + (1 | Source) +
## (1 | EngForm2) + NormHistoricalDistance + I(NormHistoricalDistance^2) +
## I(NormHistoricalDistance^3)
## Data: d2
## Control: glmerControl(optimizer = "bobyqa")
##
##      AIC      BIC   logLik deviance df.resid
## 21337.5 21393.1 -10657.8  21315.5     1150
##
## Scaled residuals:
##      Min       1Q   Median       3Q      Max
## -103.859   -1.232    0.131    1.154   171.132
##
## Random effects:
## Groups   Name                Variance Std.Dev.
## EngForm2 (Intercept)  4.124      2.031
## Set      (Intercept)  5.101      2.259
## Source   (Intercept)  4.403      2.098
## Number of obs: 1161, groups: EngForm2, 128; Set, 67; Source, 13
##
## Fixed effects:
##
##              Estimate Std. Error z value Pr(>|z|)
```

```
## (Intercept)          0.66706      0.82054      0.813      0.41624
## Age1                 1.05926      1.27451      0.831      0.40591
## Age2                -1.25631      1.63562     -0.768      0.44243
## Age3                -0.01351      1.12230     -0.012      0.99039
## AlliterationTRUE      2.28058      0.07127     31.998 < 2e-16 ***
## NormHistoricalDistance -7.46651      1.45307     -5.138 2.77e-07 ***
## I(NormHistoricalDistance^2) 10.76780      3.80645      2.829 0.00467 **
## I(NormHistoricalDistance^3) -5.41228      2.68322     -2.017 0.04369 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Correlation of Fixed Effects:
##      (Intr) Age1   Age2   Age3   AlTRUE NrmHsD I(NHD^2
## Age1      0.010
## Age2      0.407 -0.517
## Age3     -0.188 -0.263 -0.488
## AlltrtnTRUE -0.014  0.015 -0.010  0.000
## NrmHstrclDs -0.132 -0.003 -0.001  0.006 -0.029
## I(NrmHsD^2)  0.119  0.003  0.001 -0.006  0.036 -0.978
## I(NrmHsD^3) -0.110 -0.003 -0.001  0.006 -0.046  0.934 -0.985
```

```
plot_model(m3,'eff', terms="Age [all]")
```



```
pNormHist =
  plot_model(m3,'eff', terms="NormHistoricalDistance[all]") +
    xlab("Historical Distance\nBetween Forms") +
    ylab("Proportion of Norse Terms") +
    ggtitle("") +
    coord_cartesian(ylim=c(0,1))
pNormHist
```



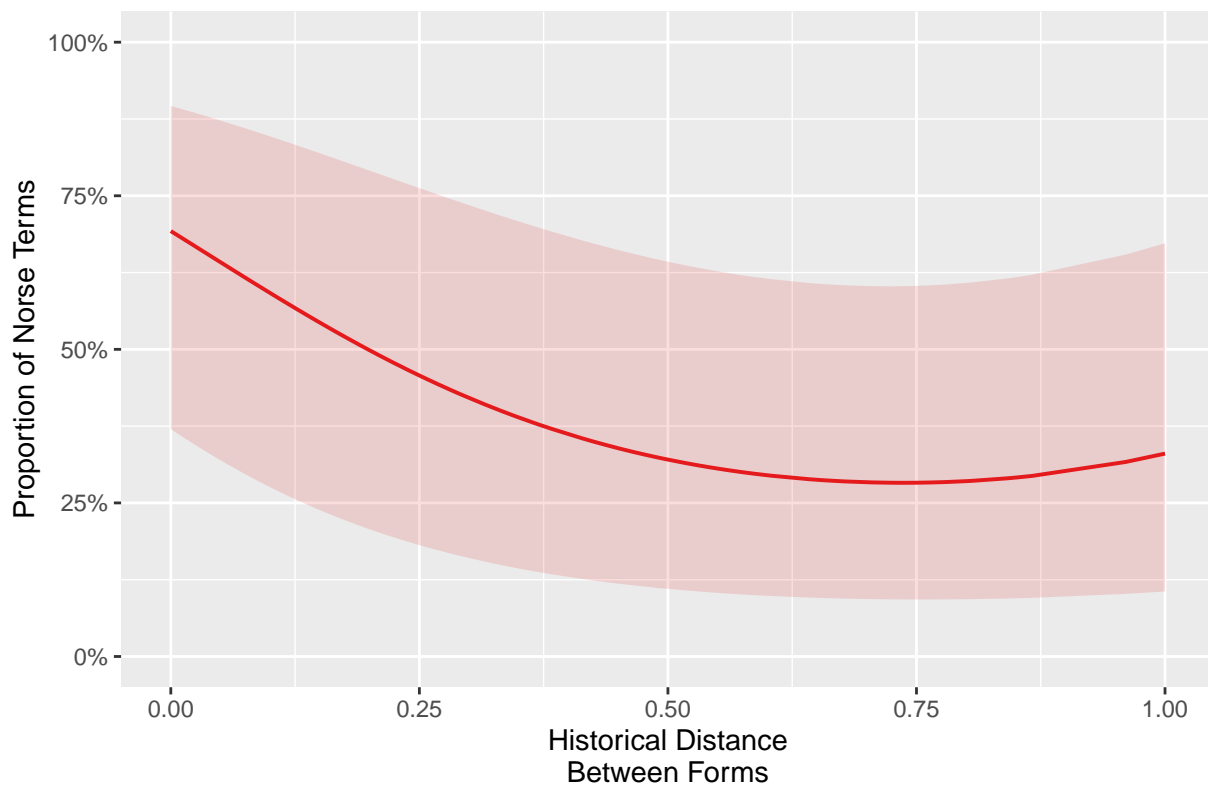
And for comparison, the quadratic model:

```
summary(m2)
```

```
## Generalized linear mixed model fit by maximum likelihood (Laplace
## Approximation) [glmerMod]
## Family: binomial ( logit )
## Formula: cbind(NFreq, EFreq) ~ Age + Alliteration + (1 | Set) + (1 | Source) +
## (1 | EngForm2) + NormHistoricalDistance + I(NormHistoricalDistance^2)
## Data: d2
## Control: glmerControl(optimizer = "bobyqa")
##
##      AIC      BIC   logLik deviance df.resid
## 21339.2 21389.7 -10659.6 21319.2    1151
##
## Scaled residuals:
##      Min       1Q   Median       3Q      Max
## -103.782   -1.231    0.131    1.158   171.091
##
## Random effects:
##  Groups   Name                Variance Std.Dev.
## EngForm2 (Intercept) 4.112      2.028
## Set       (Intercept) 5.084      2.255
## Source   (Intercept) 4.398      2.097
## Number of obs: 1161, groups: EngForm2, 128; Set, 67; Source, 13
##
## Fixed effects:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    0.48523   0.81500   0.595    0.552
## Age1           1.05582   1.28019   0.825    0.410
## Age2          -1.25146   1.64777  -0.759    0.448
## Age3          -0.01356   1.12724  -0.012    0.990
## AlliterationTRUE 2.27432   0.07118 31.950 < 2e-16 ***
## NormHistoricalDistance -4.73430   0.52341 -9.045 < 2e-16 ***
```

```
## I(NormHistoricalDistance^2) 3.21561 0.66414 4.842 1.29e-06 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Correlation of Fixed Effects:
##          (Intr) Age1  Age2  Age3  AlTRUE NrmHsD
## Age1      0.008
## Age2      0.410 -0.519
## Age3     -0.190 -0.261 -0.491
## AlltrtnTRUE -0.019 0.015 -0.010 0.000
## NrmHstrclDs -0.084 0.001 -0.002 0.001 0.040
## I(NrmHsD^2) 0.065 -0.001 0.003 -0.001 -0.052 -0.940
```

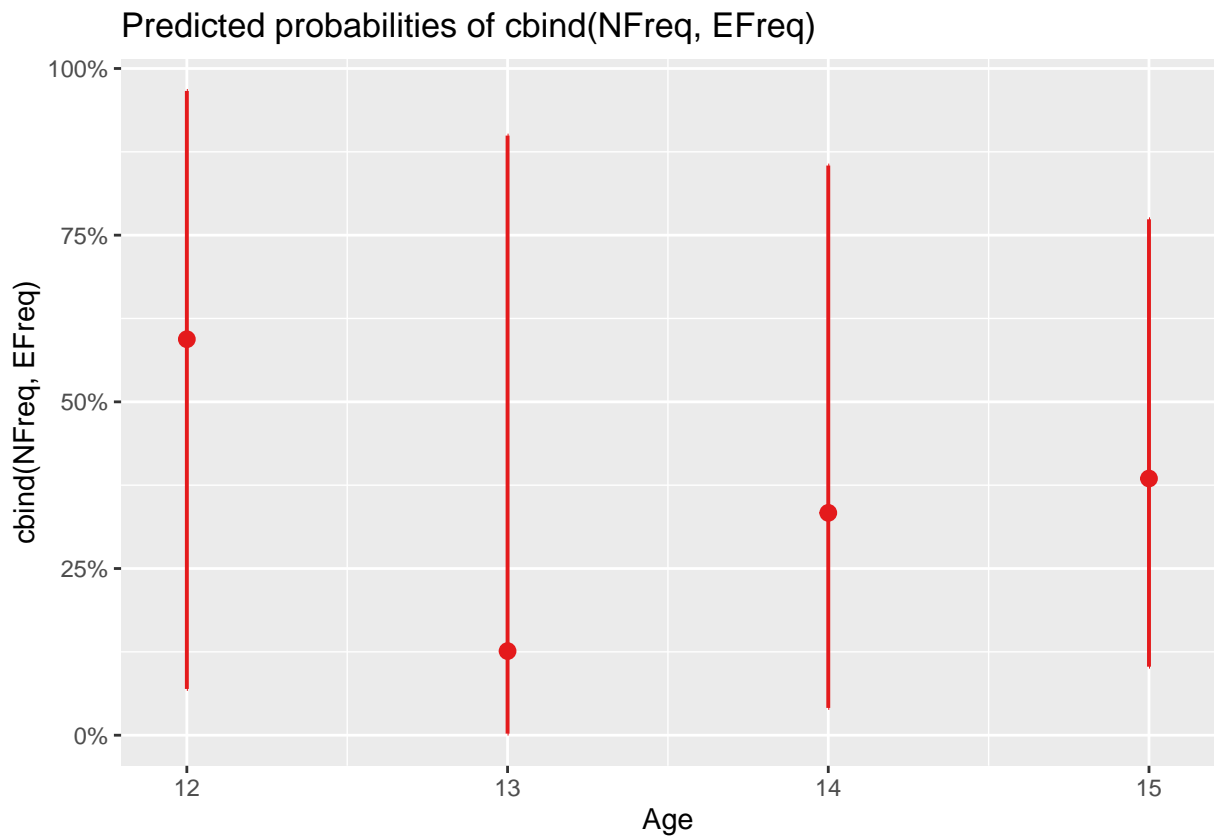
```
plot_model(m2, 'eff', terms="NormHistoricalDistance[all]") +
  xlab("Historical Distance\nBetween Forms") +
  ylab("Proportion of Norse Terms") +
  ggtitle("") +
  coord_cartesian(ylim=c(0,1))
```



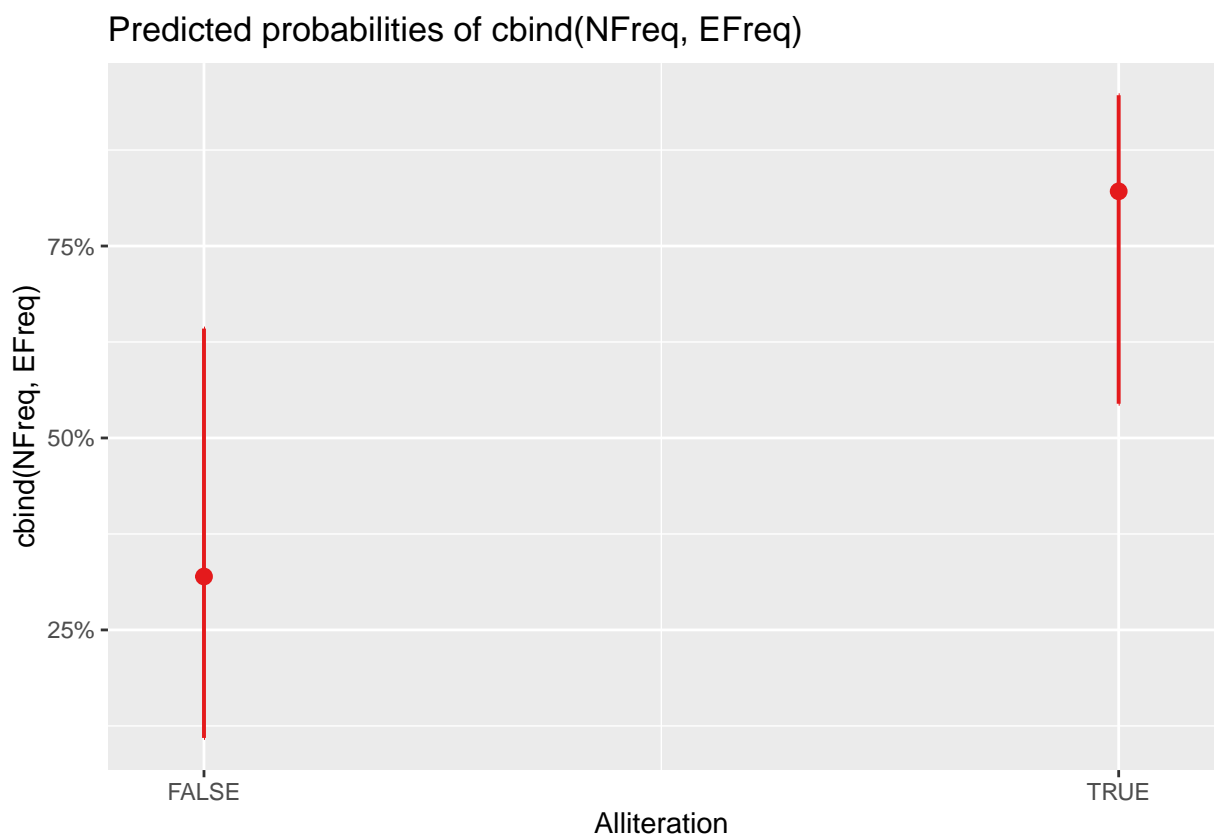
There is little qualitative difference, so for sake of easy comparison with the other results, we use the cubic model.

Effects over time and of alliteration:

```
plot_model(m3, 'eff', terms="Age [all]")
```



```
plot_model(m3, 'eff', terms="Alliteration")
```



GAM model:

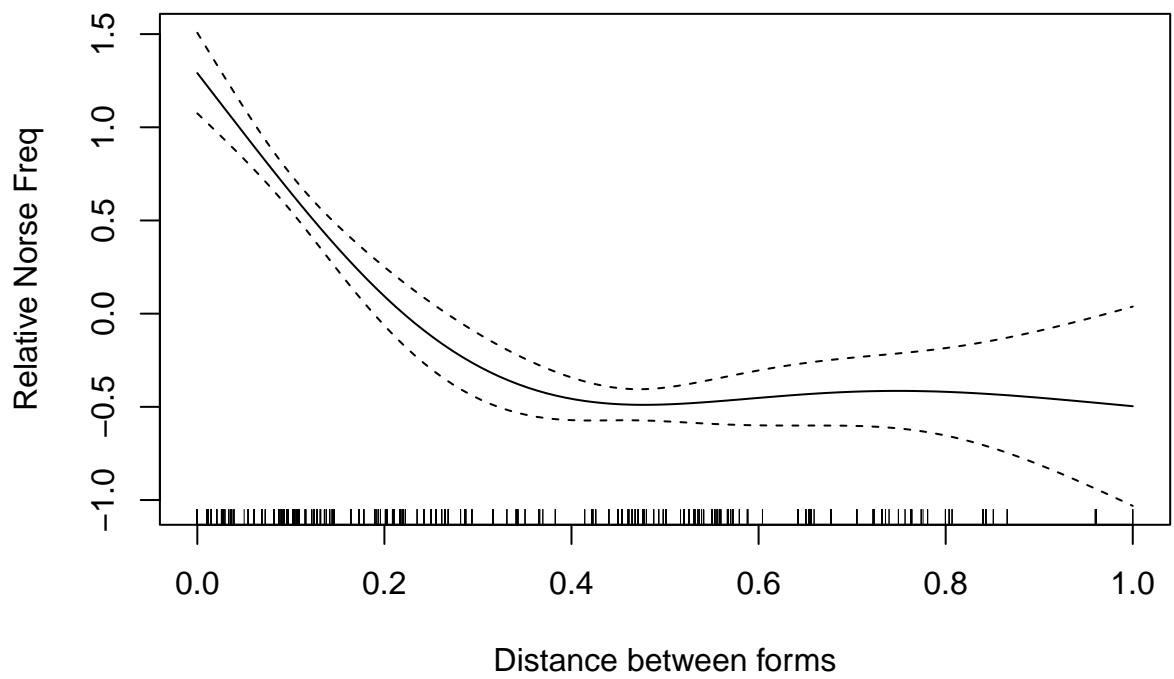
```
mGAM = gam(cbind(NFreq, EFreq) ~ Alliteration +  
            s(NormHistoricalDistance, k=4) +
```

```

      s(Age,bs = "re") +
      s(Set,bs="re") +
      s(Source,bs="re")+
      s(EngForm2,bs="re"),
      data = d2, family = "binomial")
summary(mGAM)

##
## Family: binomial
## Link function: logit
##
## Formula:
## cbind(NFreq, EFreq) ~ Alliteration + s(NormHistoricalDistance,
##   k = 4) + s(Age, bs = "re") + s(Set, bs = "re") + s(Source,
##   bs = "re") + s(EngForm2, bs = "re")
##
## Parametric coefficients:
##               Estimate Std. Error z value Pr(>|z|)
## (Intercept)   -0.42177    2.29511  -0.184    0.854
## AlliterationTRUE 2.27080    0.07131  31.843   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Approximate significance of smooth terms:
##               edf   Ref.df    Chi.sq  p-value
## s(NormHistoricalDistance) 2.918e+00   2.992 2.150e+02 < 2e-16 ***
## s(Age)                    9.446e-04   3.000 7.990e-01 0.999258
## s(Set)                    4.792e+01  66.000 5.630e+07 0.000913 ***
## s(Source)                 1.200e+01  12.000 1.147e+07 < 2e-16 ***
## s(EngForm2)               7.203e+01 127.000 2.248e+07 0.002078 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## R-sq.(adj) = 0.726   Deviance explained = 70.4%
## UBRE = 16.473   Scale est. = 1           n = 1161
plot.gam(mGAM,
  ylab="Relative Norse Freq",
  xlab="Distance between forms",
  select = 1)

```



Linear: $\beta = -7.467$, $z = -5.138$, Wald $p < 0.001$, LLDiff = 95, $df = 1$, $p < 0.001$; Quadratic: $\beta = 10.768$, $z = 2.829$, Wald $p = 0.005$, LLDiff = 11.8, $df = 1$, $p < 0.001$; Cubic: $\beta = -5.412$, $z = -2.017$, Wald $p = 0.044$, LLDiff = 1.8, $df = 1$, $p = 0.056$


```

bigPlot = grid.arrange(
  pNormSimple +
    ggtitle("Source-Level Analysis"),
  pNormFeature+
    ggtitle("") +
    theme(axis.title.y=element_blank(),
          axis.text.y=element_blank()),
  pNormHist+
    ggtitle("") +
    theme(axis.title.y=element_blank(),
          axis.text.y=element_blank()),
  nrow=1,widths=c(1.3,1,1))

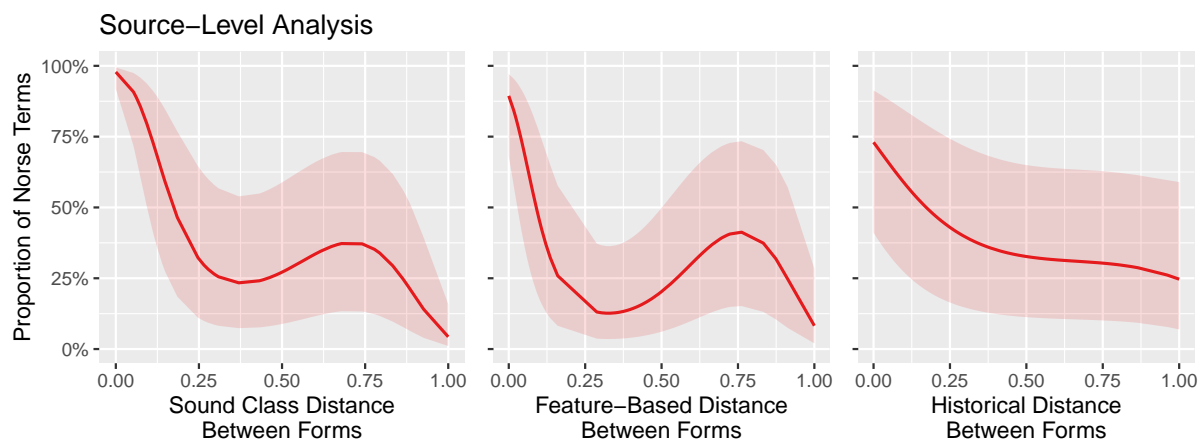
pdf("../results/BigEffectsPlot.pdf",width = 8,height=3)
plot(bigPlot)
dev.off()

```

```

## pdf
## 2

```



```

tSimple = read.csv("../results/SimpleRes_totalFreq.csv",
  stringsAsFactors = F)
tFeature = read.csv("../results/FeatureRes_totalFreq.csv",
  stringsAsFactors = F)
tHistorical = read.csv("../results/HistoricalRes_totalFreq.csv",
  stringsAsFactors = F)
sSimple = read.csv("../results/SimpleRes.csv",
  stringsAsFactors = F)
sFeature = read.csv("../results/FeatureRes.csv",
  stringsAsFactors = F)
sHistorical = read.csv("../results/HistoricalRes.csv",
  stringsAsFactors = F)

res = rbind(tSimple,tFeature,tHistorical)
resNames = c(lldiff="Log Likelihood Difference",
  Chisq="Chi Squared",pChi="p")
resNames2 = names(res)
resNames2[resNames2 %in% names(resNames)] =
  resNames[resNames2[resNames2 %in% names(resNames)]]

res = res[,names(res)!="X"]
res = cbind(Measure=rep(c("Sound Class","Feature","Historical"),each=3), res)
res = rbind(c("", "", "Model Comparison", "", "", "", "",
  "Model Estimate", "", "")),

```

```

        resNames2,
        res)

write.table(res,file="../results/MainResults_totalFreq.csv",
            sep = ",",col.names = F,row.names = F,fileEncoding = "UTF-8")

res2 = res
res2[3:11,] = rbind(sSimple,sFeature,sHistorical)

write.table(res2,file="../results/MainResults_sourceLevel.csv",
            sep = ",",col.names = F,row.names = F,fileEncoding = "UTF-8")

```

7 Comparison between texts

Compare sources according to the difference in proportion of Norse terms for each set.

```
g = read.xlsx("../data/SharedIntegrationOfCognatesData.xlsx",1)
g = g[g$Etymology %in% c("Norse","English"),]
# Ignore numerals in set name
g$Set = gsub("[0-9]", "", g$Set)

nCols = names(g)[which(names(g)=="No..in.Ormulum"):which(names(g)=="Rolle")]
for(col in nCols){
  g[,col] = as.numeric(g[,col])
}

## Warning: NAs introduced by coercion

## Warning: NAs introduced by coercion

## Warning: NAs introduced by coercion

allSets = unique(g$Set)
f = data.frame(Set = allSets)
fLog = data.frame(Set = allSets)
for(col in nCols){
  gN = g[g$Etymology == "Norse",]
  gE = g[g$Etymology == "English",]
  fNorse = tapply(gN[,col], gN$Set, sum)[allSets]
  fNorse[is.na(fNorse)] = 0
  fEng = tapply(gE[,col], gE$Set, sum)[allSets]
  fEng[is.na(fEng)] = 0
  f[,col] = fNorse / (fEng+fNorse)
  fLog[,col] = log10(1+fEng) - log10(1+fNorse)
}

nColsLabels = c("Ormulum", "FCPC", "Havelok", "Genesis & Exodus",
                 "Mannyng", "Gawain-poet", "Wars of Alexander",
                 "St Erkenwald", "Cursor Mundi", "Lincolnshire",
                 "Nottinghamshire", "Norfolk", "Rolle")

mat = matrix(NA, nrow=length(nCols),
             ncol = length(nCols))
rownames(mat) = nColsLabels
colnames(mat) = nColsLabels
matLog = matrix(NA, nrow=length(nCols),
               ncol = length(nCols))
rownames(matLog) = nColsLabels
colnames(matLog) = nColsLabels
for(i in 1:length(nCols)){
  iProp = f[,nCols[i]]
  iPropLog = fLog[,nCols[i]]
  for(j in 1:length(nCols)){
    jProp = f[,nCols[j]]
    jPropLog = fLog[,nCols[j]]

    diffs = abs(iProp-jProp)
    diffs = diffs[!is.nan(diffs)]
    diffs = diffs[!is.na(diffs)]
    mat[nColsLabels[i], nColsLabels[j]] = mean(diffs)
  }
}
```

```

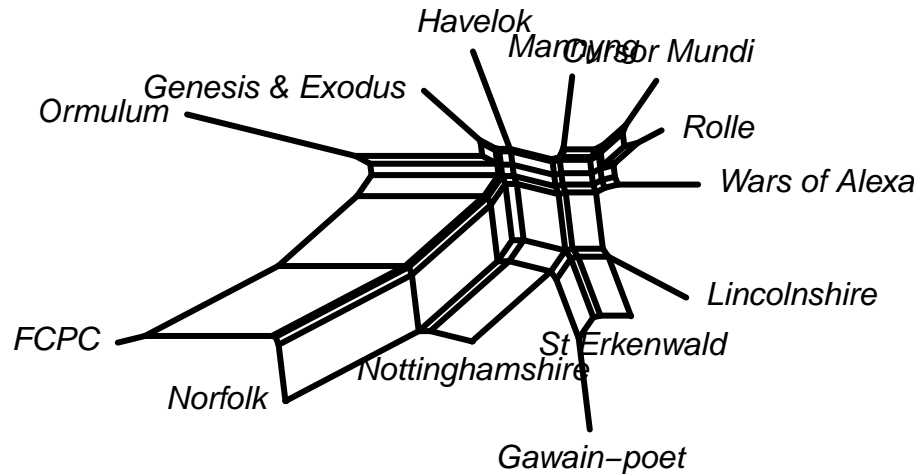
    diffsLog = abs(iPropLog - jPropLog)
    matLog[nColsLabels[i],nColsLabels[j]] = mean(diffsLog)
  }
}

```

```

phy = neighborNet(mat)
plot(phy)

```



```

pdf("../results/NeighbourNet.pdf")
plot(phy)
dev.off()

```

```

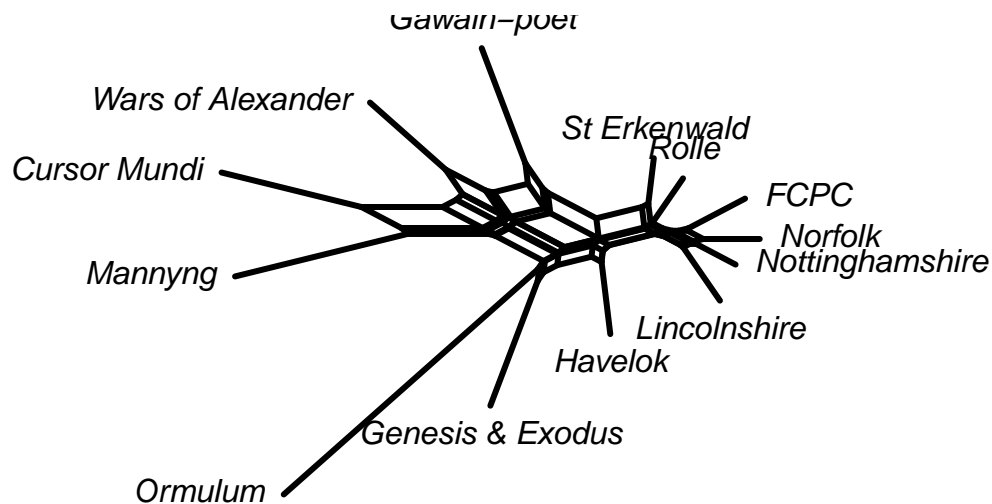
## pdf
## 2

```

```

phyLog = neighborNet(matLog)
plot(phyLog)

```



```

pdf("../results/NeighbourNet_Log.pdf")
plot(phyLog)
dev.off()

```

```

## pdf
## 2

```