

This sequence picks a single height value per Gen2 subject.

Define the age cutoffs to keep ages within the same Window as Gen1 Heights. Define the height cutoffs to exclude values that are more likely to be entry errors or a developmental disorder, than a true reflection of additive genetics

Load the appropriate information from the SQL Server database

```
summary(dsSubjectYear)
```

```
comma(c(nrow(dsHeightLong), nrow(dsSubjectYear)))
```

```
[1] "70,614" "86,579"
```

Combine the feet and inches to get total inches. Filter out records with height values that are outside the desired range

[1] 35307

[1] 35067

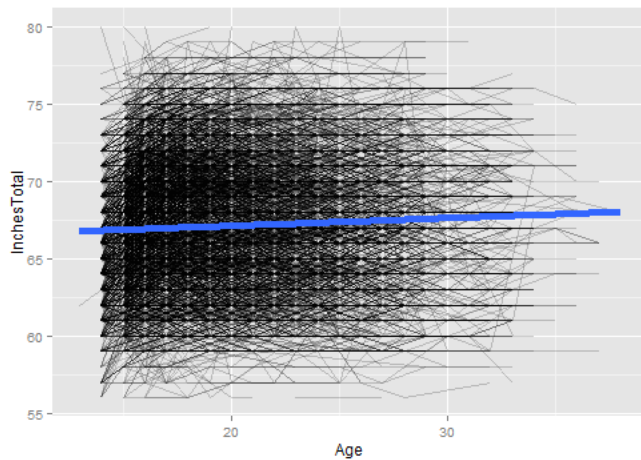
SubjectTag	SurveyYear	InchesTotal
Min. : 301	Min. :1994	Min. :56.0
1st Qu.: 267502	1st Qu.:2002	1st Qu.:64.0
Median : 546701	Median :2006	Median :67.0
Mean : 549932	Mean :2004	Mean :67.2
3rd Qu.: 805901	3rd Qu.:2008	3rd Qu.:70.0
Max. :1266703	Max. :2010	Max. :80.0

A histogram showing the distribution of heights (InchesTotal) for the variable dsHeightYear. The x-axis is labeled 'dsHeightYear\$InchesTotal' and ranges from 55 to 80. The y-axis is labeled 'count' and ranges from 0 to 3000. The distribution is roughly bell-shaped, centered around 66 inches, with a peak count of approximately 3200.

Join the height data with age of the subject when the height was taken. Filter out records where the age is outside of the desired window.

[1] 35067

QR



```
rm(dsSubjectYear, dsHeightYear)
```

Standardize by Gender & Age. Calculated Age (using SurveyDate and MOB) has been truncated to integers.

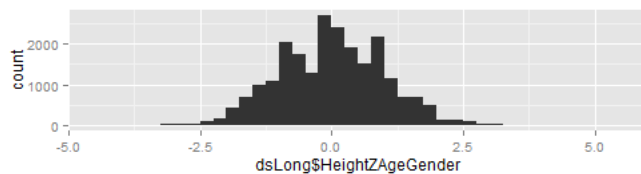
```
dsLong <- dsLong[ageMin <= dsLong$Age & dsLong$Age <= ageMax, ]
nrow(dsLong)
```

```
[1] 22795
```

```
dsLong <- ddpoly(dsLong, c("Age", "Gender"), transform, HeightZAgeGender=scale(InchesTotal))
nrow(dsLong)
```

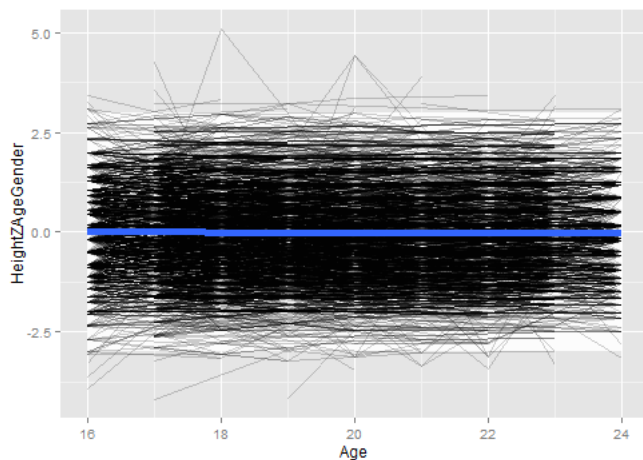
```
[1] 22795
```

```
qplot(dsLong$HeightZAgeGender, binwidth=.25) #Make sure ages are normalish with no extreme values.
```



Determine Z-score to clip at. Adjust as necessary (zMin & zMax were defined at the top of the page). The white box extends between zMin and zMax.

```
ggplot(dsLong, aes(x=Age, y=HeightZAgeGender, group=SubjectTag)) +
  annotate("rect", xmin=min(dsLong$Age), xmax=max(dsLong$Age), ymin=zMin, ymax= zMax, fill="gray99") +
  geom_line(alpha=.2) + geom_smooth(method="rlm", aes(group=NA), size=2)
```



```
dsLong <- dsLong[zmin <= dsLong$HeightZAgeGender & dsLong$HeightZAgeGender <= zMax, ]
nrow(dsLong)
```

```
[1] 22733
```

Pick the subject's oldest record (within that age window). Then examine the age & Z values

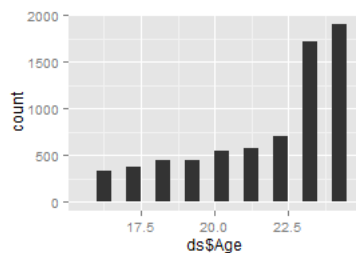
```
ds <- ddpoly(dsLong, "SubjectTag", subset, rank(-Age)==1)
summary(ds)
```

SubjectTag	SurveyYear	Age	Generation	Gender
Min. : 301	Min. :1994	Min. :16.0	Min. :2	Min. :1.00
1st Qu.: 266202	1st Qu.:2004	1st Qu.:20.0	1st Qu.:2	1st Qu.:1.00
Median : 537401	Median :2008	Median :23.0	Median :2	Median :1.00
Mean : 545706	Mean :2007	Mean :21.5	Mean :2	Mean :1.49
3rd Qu.: 804403	3rd Qu.:2010	3rd Qu.:24.0	3rd Qu.:2	3rd Qu.:2.00
Max. :1266703	Max. :2010	Max. :24.0	Max. :2	Max. :2.00
InchesTotal	HeightZAgeGender			
Min. :56.0	Min. :-2.9855			
1st Qu.:64.0	1st Qu.: -0.7195			
Median :67.0	Median : -0.0730			
Mean :67.5	Mean : -0.0016			
3rd Qu.:71.0	3rd Qu.: 0.5766			
Max. :79.0	Max. : 2.9905			

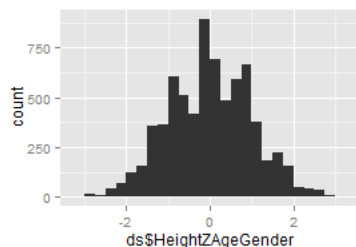
```
# SELECT [Mob], [LastSurveyYearCompleted], [AgeAtLastSurvey]
# FROM [NlsLinks].[dbo].[vewSubjectDetails79]
# WHERE Generation=2 and AgeAtLastSurvey >=16
#After the 2010 survey, there were 7,201 subjects who were at least 16 at the last survey.
nrow(ds)
```

```
[1] 7069
```

```
qplot(ds$Age, binwidth=.5) #Make sure ages are within window, and favoring older values
```



```
qplot(ds$HeightZAgeGender, binwidth=.25) #Make sure ages are normalish with no extreme values.
```



Compare with Kelly's height values. Make sure they roughly agree. There are a few differences, including (1) the age range is a little shifted, (2) the 2010 survey wasn't available, (3) the cutoff scores were more generous, and (4) the order of standardization & clipping *might* have been different.

```
# Compare against kelly's previous versions of Gen2 Height
# pathInputkellyOutcomes <- "F:/Projects/Nls/Links2011/CodingUtilities/Gen2Height/ExtraOutcomes79FromKelly"
dskelly <- read.csv(pathInputkellyOutcomes, stringsAsFactors=FALSE)
dskelly <- dskelly[, c("SubjectTag", "HeightStandardizedFor19to25")]
dsOldvsNew <- join(x=ds, y=dskelly, by="SubjectTag", type="full")

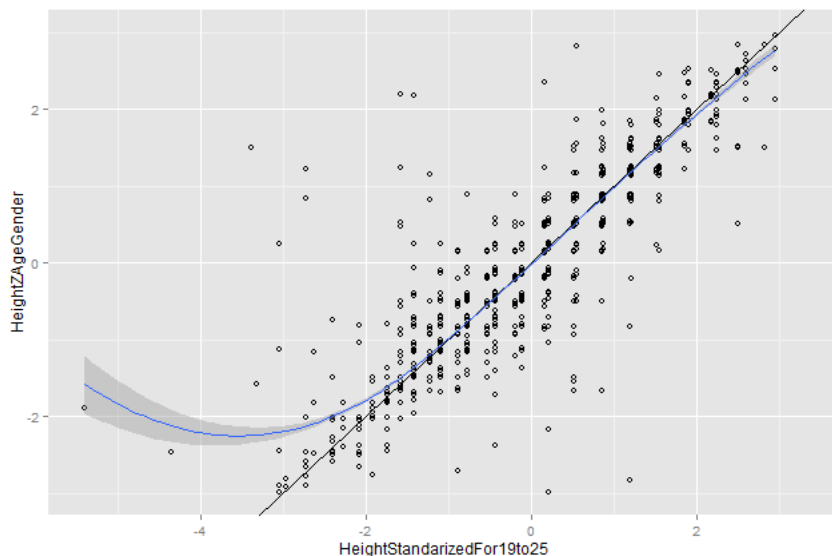
#See if the new version is missing a lot of values that the old version caught.
# The denominator isn't exactly right, because it doesn't account for the 2010 values missing in the new
table(is.na(dsOldvsNew$HeightzAgeGender), is.na(dsOldvsNew$HeightStandardizedFor19to25), dnn=c("NewIsMissing", "OldIsMissing"))
```

	oldIsMissing	
NewIsMissing	FALSE	TRUE
FALSE	5089	1980
TRUE	34	4392

```
#View the correlation
cor(dsOldvsNew$HeightZAgeGender, dsOldvsNew$HeightStandardizedFor19to25, use="complete.obs")
```

```
[1] 0.9553
```

```
#Compare against an x=y identity line.
ggplot(dsOldvsNew, aes(x=HeightStandardizedFor19to25, y=HeightZAgeGender)) + geom_point(shape=1) + geom_a
```



```
# @knitr writeToCsv
write.csv(ds, pathOutputSubjectHeight, row.names=FALSE)
```