

Class 04 Notes

Isabella Chittumuri

9/23/2020

```
# How to get working directory
getwd()

## [1] "/Users/isabellachittumuri/Documents/Hunter College/Fall 2020/Stat 706/Class notes"
knitr::opts_chunk$set(echo = TRUE)

library(faraway)
library(dplyr)

##
## Attaching package: 'dplyr'
## The following objects are masked from 'package:stats':
##
##   filter, lag
## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union

library(ggplot2)
library(gridExtra)

##
## Attaching package: 'gridExtra'
## The following object is masked from 'package:dplyr':
##
##   combine

library(printr)

## Registered S3 method overwritten by 'printr':
##   method          from
##   knit_print.data.frame rmarkdown

theme_set(theme_minimal()) # automatically set a simpler ggplot2 theme for all graphics
```

Binary Data

The example data

- Heart Disease Data Example, in the wgs data set
- Chronic heart disease (chd) is going to be the response, what we're trying to predict/understand. This is a binary variable, only takes two values. In this case, it's a character vector, meaning the values are

- yes and no.
- Other variables, height in inches and cigs per day, are both continuous variables.

```
summary(wcgs[, c("chd", "height", "cigs")])
```

chd	height	cigs
no :2897	Min. :60.00	Min. : 0.0
yes: 257	1st Qu.:68.00	1st Qu.: 0.0
NA	Median :70.00	Median : 0.0
NA	Mean :69.78	Mean :11.6
NA	3rd Qu.:72.00	3rd Qu.:20.0
NA	Max. :78.00	Max. :99.0

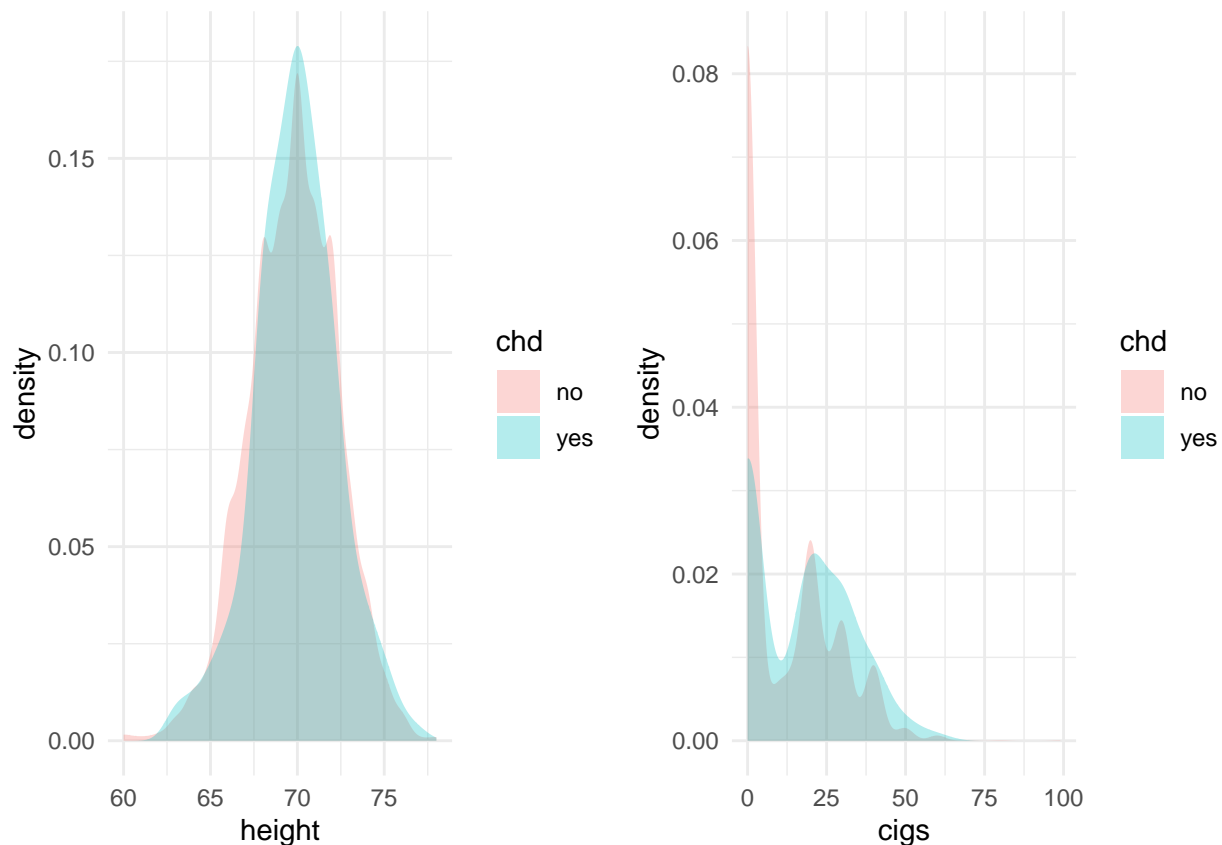
- Something that stands out is the max number of cigs a day, 99 per day
- The median is zero, meaning 50% of people don't smoke

Example Data | Plot

```
# Density plots are like histograms scaled to 1 so we can compare
p1 <- ggplot(wcgs, aes(fill = chd, x = height)) +
  geom_density(alpha = .3, color = NA)

p2 <- ggplot(wcgs, aes(fill = chd, x = cigs)) +
  geom_density(alpha = .3, color = NA)

grid.arrange(p1, p2, ncol = 2)
```



- First: there is a lot of overlap between the height of people and the density of chd. This suggests that there probably isn't a lot of relationship between those two values.
- Second: you can see a pattern that the red is way higher on the left. From what we know, if you smoke more than 0 packs a day there might be a relationship between chd.
- There are little peaks and drops. One in particular is the big drop which is around 20 cigs which is a pack of cigs a day. Seems like not a lot of ppl smoke 0-20 cigs, but once you get to 20 more a lot more people start to smoke

Let's create a linear model

```
# The %>% (pip) notation is a way of putting the wcgs mutate(wcgs2, chd). We're basically passing wcgs2
wcgs2 <- wcgs %>%
  mutate(chd = as.numeric(chd == 'yes'))

lm(chd ~ height, data = wcgs2)
```

```
##
## Call:
## lm(formula = chd ~ height, data = wcgs2)
##
## Coefficients:
## (Intercept)      height
##    -0.06083      0.00204
```

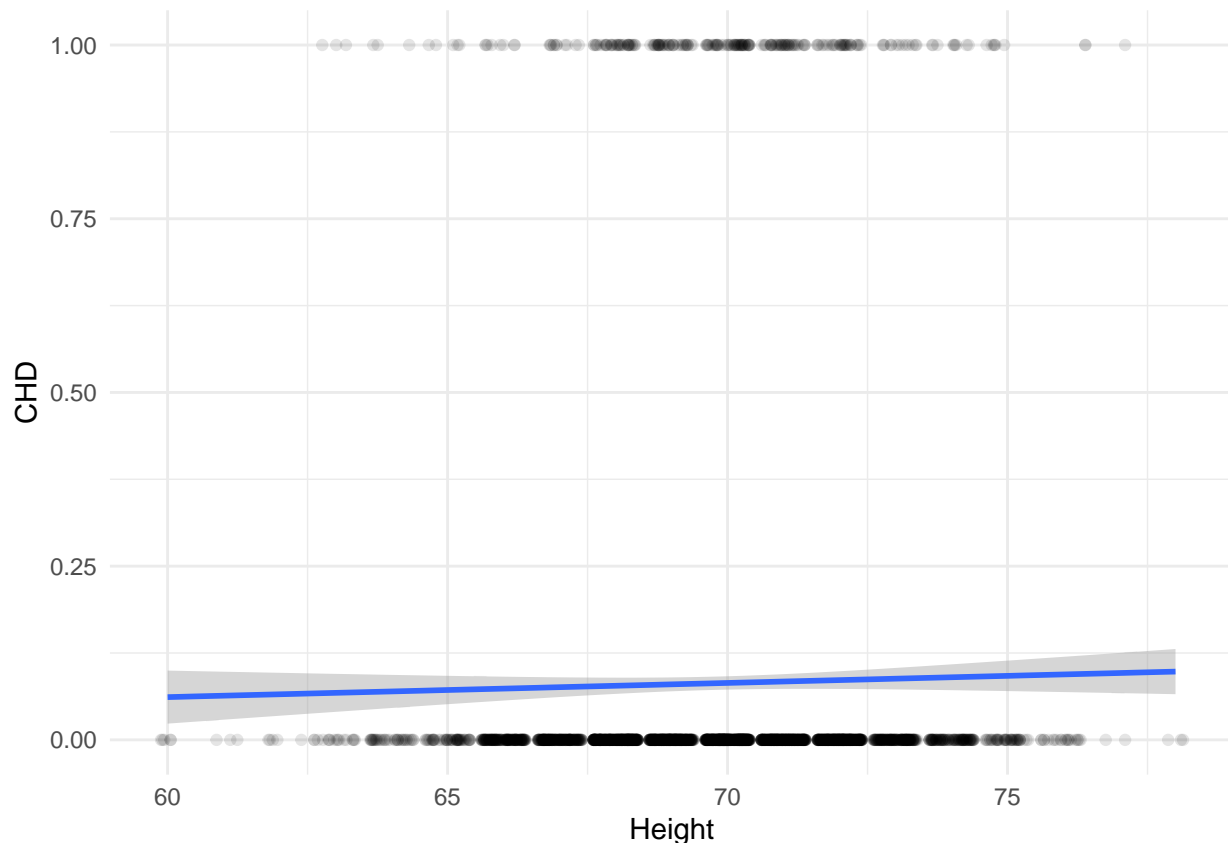
- We created a copy of the data set and turned the yes/no field into 1/0. This becomes a true/false and you turn it into a number
- Generally, the thing you're interested in, in this case people who do have chd, you typically turn that into a 1.

Let's create a linear model

- Smooth function is similar to confidence interval

```
#scatter plot between height adn chd
ggplot(wcgs2, aes(x = height, y = chd)) +
  geom_jitter(height = 0, alpha = 0.1) +
  labs(x = "Height", y = "CHD") +
  geom_smooth(method = "lm")
```

```
## `geom_smooth()` using formula 'y ~ x'
```



- The blue line is the linear model we created
- The dots are the actual data points, so when you have chd as yes it becomes a 1 and when chd is a no it becomes a 0.

What's the issue with this model? - Note: when height is zero, our output is a negative number, the intercept -0.06. We want a probability number between 0-1. This model can only be 0 or 1, but we want a probability. The model doesn't really work correctly because it can give us negative numbers. Similarly, if we have a height of 1000, our prediction for chd is going to be above 1 which is non seneschal - You can sort of clip it by saying anything below 0 is 0, but it's not as satisfying. - This isn't a good model so we need to do something differently

Binary Regression:

- Assume response is distributed Bernoulli $B(p_i)$ for each i observation
- The Bernoulli distribution is always for one case, like a coin flip

The density function for the Bernoulli distribution below:

$$P(Y_i = y_i) = p_i^{y_i} (1 - p_i)^{1-y_i}$$

- This says that if Y is 1, the right side defines how likely something is to happen given a probability
- We want to model p, for a person with certain predictor/characteristics what is the probability that they have chd. We're going to model this value as a linear combination of different predictors
- Let's construct a linear predictor with q terms, this is how we're going to combine our other predictors into one single value (eta)
- We want to combine because we want our predictors to influence and play a part in p

- We're going to make $p = \eta$

$$\eta_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_q x_{iq}$$

- Before we had the μ as our predictor in the normal density equation.
- But we still have prob, there is no guarantee that this linear combination will fall between 0 and 1. And p in this equation has to fall between 0 and 1.

What should we use now? - We're forced to transform η into something that falls between 0 and 1. Once we transform it, we're going to stick that value for p and that's how we got the coefficients that we're looking at into the Bernoulli equation

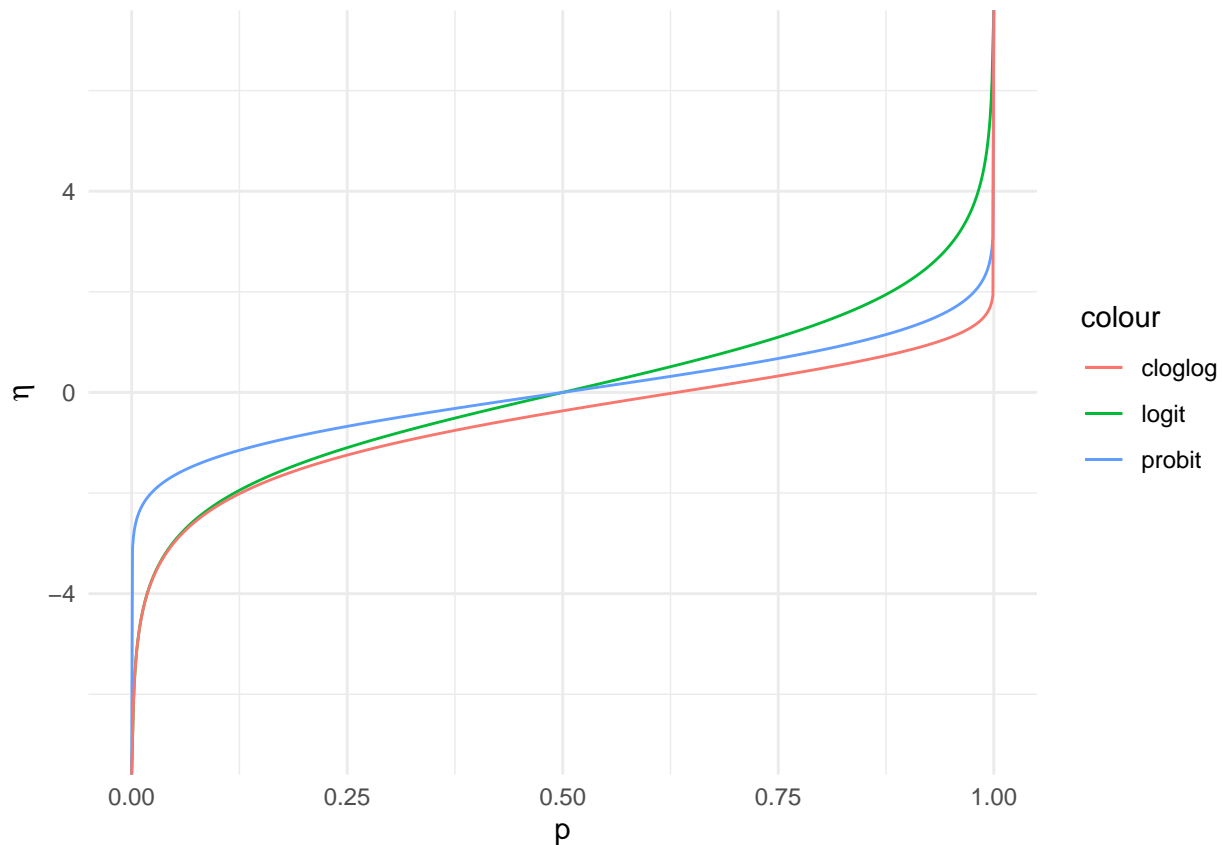
We need a link function

- A link function is often used to transform the variable of interest (in this case p_i) to a reasonable scale for the linear predictor
- Making all values from negative infinity and infinity will be transformed into a value between 0 and 1
- They are called link function because they link the parameters to the linear predictor, η is also called the linear predictor
- Options
 - Logit: $\eta = \log(p/(1-p))$ Above: p runs from 0-1, so basically you have $p/(1-p)$ run from 0 to infinity. This is also called the odds, odds run from 0 to infinity. If we take the log of that we can go from negative infinity to infinity
 - Probit: $\eta = \Phi^{-1}(p)$ where Φ is the normal cumulative distribution Above: the inverse of the normal cumulative distribution. The normal distribution can take a value between negative infinity to infinity. So we're taking the percental and changing it back into a number. How we can turn 0-1 into negative infinity to infinity
 - Complementary log-log: $\eta = \log(-\log(1-p))$ Above: does the same as second, runs the appropriate values
- We need our values go from negative infinity to infinity, because they are defined by the linear combination of η . The estimate of betas can be any number between negative infinity to infinity, like temperature. If β_{η_1} is 1 and X is temperature, the temperature can range from -100 to 100. That means that η can range from -100 to 100, prior to transformation. But we need to put it in p , but p can only go from 0-1 so that's why we are transforming it

Link function graphs

Here's an example of the three functions: logit, probit, and cloglog

```
ggplot(data.frame(p = seq(0,1, length.out = 1000)), aes(x = p)) +
  geom_line(aes(y = log(p/(1-p)), color = "logit")) +
  geom_line(aes(y = qnorm(p), color = "probit")) +
  geom_line(aes(y = log(-log(1-p)), color = "cloglog")) +
  labs(y = expression(eta))
```



- X-axis is the p , and on the y-axis is the transformed data
- Left tail asymptotically goes to 0 and right tail asymptotically goes to 1

What tools do we use to solve this?

- MLE! maximum likelihood estimation
- We're going to use the logistic regression
- Last time there were formulas for everything, but that's not going to be possible here.
- Example with logistic function

If we invert the logistic function (logit), we will get this:

$$p = \frac{e^\eta}{e^\eta + 1}$$

$$l(\beta) = \sum_{i=1}^n [y_i \eta_i - \log(1 + e_i^\eta)] \eta = \log(p/(1-p))$$

- See Appendix A for a bit more details. It used the binomial distribution but the Bernoulli is a special case of the binomial distribution so you should be able to follow
- Remember that the loglikelihood is when you can plug in
- Using the density distribution and the link function logit you can rewrite the loglikelihood
- The loglikelihood for this when using the logit function is equal to the first formula above

Run it in R

Before we ran the linear model function (lm), now we're going to use the generalized linear models (glm) function. - We're going to do it the exact same way with a small difference. - To create our formula, chd is

the response variable and the predictor is height - Now we're going to specify the family that is used and specific the link function. This means that we're going to use the binomial and logit link function. You don't need to write that every time because that link function is the default for binomial link. And then we're going to be using the wgs data - Recall, wgs had chd was a character vector, yes and no. R decides what the 1 and 0 is because when in doubt it does things alphabetically. So behind the scenes. R converts chd into a 0 and 1 variable. We would have the same exact output if we use the wgs2 data set. - Happy accident in the English language where yes comes after no in the alphabet. so you kind of auto get what you want, yes into 1 and no into 0

```
logitmod <- glm(chd ~ height, family = binomial(link = "logit"), data = wgs)
summary(logitmod)
```

```
##
## Call:
## glm(formula = chd ~ height, family = binomial(link = "logit"),
##      data = wgs)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -0.4587  -0.4186  -0.4131  -0.4024   2.3157
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -4.33732    1.81231  -2.393   0.0167 *
## height       0.02742    0.02590   1.058   0.2899
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 1781.2  on 3153  degrees of freedom
## Residual deviance: 1780.1  on 3152  degrees of freedom
## AIC: 1784.1
##
## Number of Fisher Scoring iterations: 5
```

- The fisher scoring iterations is the algorithm of how it actually solve the equation

Compare links

- We fit a probit model in this as well, it's dotted but as you can see you can't even see it in the graph. It overlaps into the logit because they are generally similar in how they act

```
probitfit <- glm(chd ~ height, family = binomial(link = "logit"), data = wgs)
```

```
fake_data <- data.frame(height = seq(10, 100, length.out = 100))
```

```
ggplot(wgs, aes(x = height, y = as.numeric(chd == 'yes')) +
```

```
  # geom_point() +
```

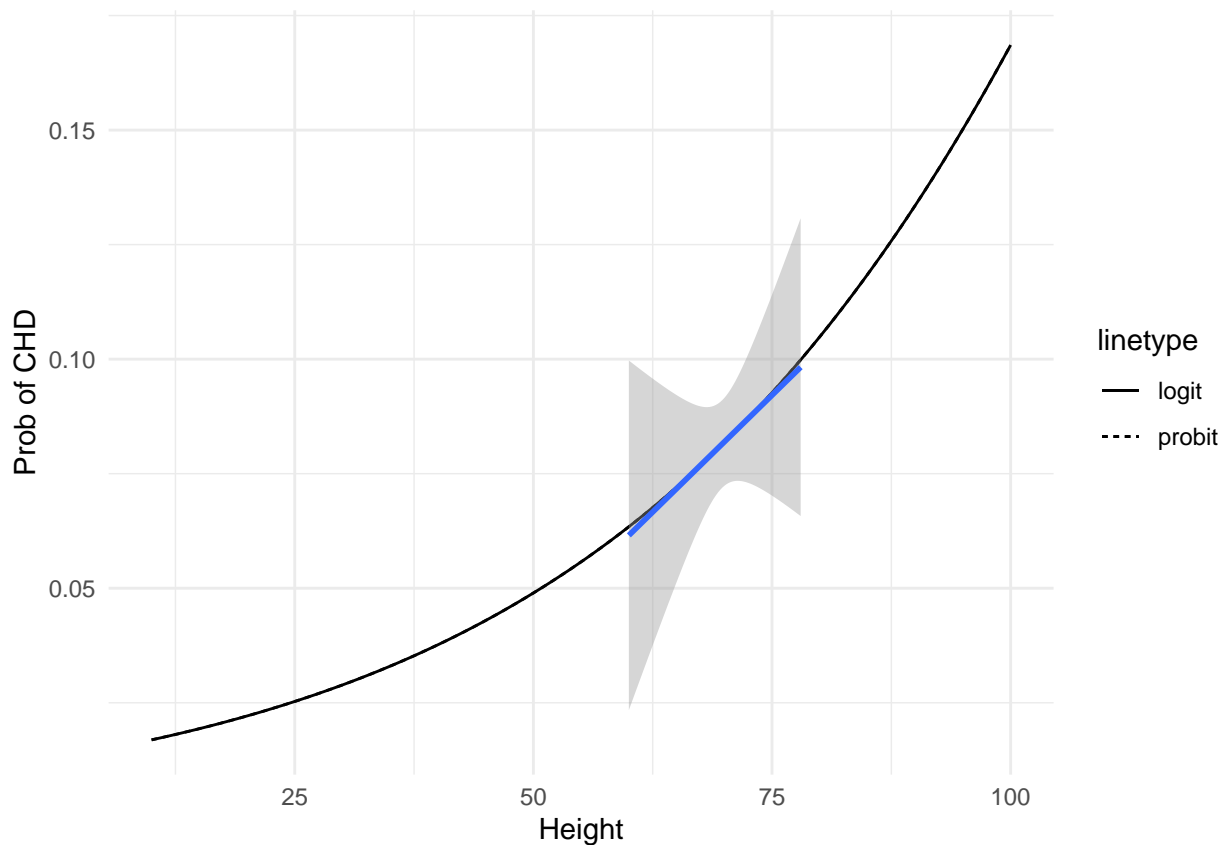
```
  labs(x = "Height", y = "Prob of CHD") +
```

```
  geom_line(data = fake_data, aes(y = predict(logitmod, fake_data, type = "response"), linetype = "logi
```

```
  geom_line(data = fake_data, aes(y = predict(probitfit, fake_data, type = "response"), linetype = "probi
```

```
  geom_smooth(method = 'lm'))
```

```
## `geom_smooth()` using formula 'y ~ x'
```



```
# coord_cartesian(ylim = c(0,1))
```

- This curved line is the estimate of prob of chd given height. It looks like it's starting to approach the probability of zero as opposed to the old line the linear version, tangent to the curve.
- Our predicted probability of chd is now curved

Fit the full model

- We're going to a larger model. So instead of one predictor of height, we're going to add cigs into the model.

```
lmod <- glm(chd ~ height + cigs, family = binomial, wgs)
summary(lmod)
```

```
##
## Call:
## glm(formula = chd ~ height + cigs, family = binomial, data = wgs)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.0041  -0.4425  -0.3630  -0.3499   2.4357
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -4.50161    1.84186  -2.444   0.0145 *
## height       0.02521    0.02633   0.957   0.3383
## cigs         0.02313    0.00404   5.724 1.04e-08 ***
## ---
```



```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 1781.2  on 3153  degrees of freedom
## Residual deviance: 1749.0  on 3151  degrees of freedom
## AIC: 1755
##
## Number of Fisher Scoring iterations: 5
```

How do we create predictions?

- We're going to be using coefficients of the model
- We want to predict the change of heart disease for someone who is 50 increases tall and who smoked 10 cigs per day.
- We create the linear predictor, eta. Remember the formula, intercept, betas and X's
- Eta can range from negative infinity to infinity so we're going to use the equation from line 142 and turn it into a probability.

```
coef(lmod)

## (Intercept)      height      cigs
## -4.50161397  0.02520779  0.02312740

# You can write little function in R to make it easy to reuse code. This is just an inverse logistic fu
# The odds is just a function, so you can stick eta into where odds is there. Odds is just a variable t
ilogit <- function(logodds){
  exp(logodds)/(exp(logodds) + 1)
}

eta <- -4.50161397 + .02520779 * 50 + 10 * 0.02312740

# Same as saying the odds
ilogit(eta)

## [1] 0.04697836

predict(lmod, data.frame(height = 50, cigs = 10), type = "response")

##      1
## 0.04697836
```

- ilogit(eta) is the prob of 4.7%
- R can do all the work by using the predict function on the linear model, put in a data frame with height of 50 and cigs of 10 and get a type response because the default is to actually show you the linear predictor and you get the same number
- The default linear predictor is that it'll give you the untransformed eta

Interpreting Odds

$$\frac{p}{1-p} = op = \frac{o}{1+o}$$

- Odds run unbounded from $(0, \infty)$ and is defined as the first equation, same type of odds for gambling and betting
- You can convert odds to p-values by transforming the odd by logging it and getting the second equation which will give you the log odds from $(-\infty, \infty)$
- You can interpret odd by odds

- And odds of greater than 1 means that p is greater than 1-p, meaning it is more likely to happen
- The odds are only 1 if p is .5 and 1-p is .5
- The odds are nice because they go from 0 to infinity instead of being bound by 1

Odds in regression

- First is both equal to eta, the odds is = eta to the b0...

$$\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 x_1 \text{ odds} = e^{\beta_0} e^{\beta_1 x_1}$$

- This only applies to logit, hard to explain with probit, cloglog

```
exp(coef(lmod))
```

```
## (Intercept)      height      cigs
## 0.01109108  1.02552819  1.02339691
```

- What you always do is report the exponentiated versions of the coefficients
- The result is the odds ratios, the e to the betas..
- From this we can say that the odds of heart disease increase by 2.6% with each additional inch in height and by 2.3% with each additional cigarette smoked per day
- The way this works is that we have e to the β_1 , we can take the x value out, so it's to the power. So that every time you go up by 1 in x, you go up another 1.02 on odds, same as raising on 2.3%
- e to a value is always going to be positive, e to a neg value is always going to be less than 1
- From this result, your odds go up by 1.02. If you use % for odds, people sometimes think you're taking about relative risk.

Odds

- let's say x_1 is a dummy variable (1 or 0)
- when $X_1 = 1$, it equal first eq
- when $X_1 = 0$, you only get e to the β_0 , second eq

$$\text{odds}_1 = e^{\beta_0} e^{\beta_1} \text{odds}_0 = e^{\beta_0}$$

- Odds ratio
- It's the ratio of those values previously

$$\frac{\text{odds}_1}{\text{odds}_0} = e^{\beta_1}$$

- This is always going to be the case because so many things cancel out
- When we hold everything at a constant we can interpret the change these effects have on the outcome individually
- The odds are mutually exclusive, it's always p/1-p
- Odds can be seen as the probability, you can find the probability from the odds and vice versa.

Relative Risk

- Relative risk is another measure that can be easier to
- Instead of having the odds of 1 divided the odds of 0, we're going to look at the probability of 1 / prob of 0
- This uses probability, we don't use probability when using odds
- Relative risk is probability of success given 1 vector/ prob of success given another vector
- Example what is the relative risk if I smoke another pack a day?

```
# Create the prob someone has chd if they smoke 20 cigs, 68"
ilogit(sum(coef(lmod) * c(1,68, 20)))
```

```
## [1] 0.08907868
```

```
# 0 cigs, 68"
```

```
ilogit(sum(coef(lmod) * c(1,68, 0)))
```

```
## [1] 0.05800425
```

- Relative risk is just the division of those two values above

```
.089/.058
```

```
## [1] 1.534483
```

- With the odds ratio, you don't have to compute everything, you can just look at individual pieces. You would look at the $\exp(\text{coef}(\text{lmod}))$ at the cigs portion, this is what 1 cigs does. If you want to say what 20 cigs does, you would do $(1.02)^{20}$ which gives you 1.48, pretty close to the relative risk
- Often times if you are taking about values that are 10% or lower, odds ratio and relative risk follow each other
- The relative risk of smoking 20 cigs relative to smoking 0 cigs is a 53% increase for developing chd
- You can get in trouble saying something and it being interpreted as relative risk when you really meant an odd

Inference

- Let's talk about other ways to actually test the model
- From the deviance slide, you can see that height doesn't actually seem to be significant, however the z-value for that is an approximation.
- But really we want to use something fancier to compare two models. We can do that using the deviance is the difference between a prefect model and a smaller model

Likelihood Ratio

???

diff between likelihood, loglikelihood, deviance, and logarithmic scoring.

- L_L is a larger (prefect) model, L_S is a smaller model

$$2\log \frac{L_L}{L_S}$$

- This is the definition of the deviance in general for all cases. What we're saying is if the large model is a prefect fitting, the loglikelihood is zero, meaning that this basically goes away. And we're left with the equation below.
- Use saturated model where L_L fits the data perfectly so $L_L = 0$.

$$D = -2 \sum_{i=1}^n \hat{p}_i \log(\hat{p}_i) + \log(1 - \hat{p}_i) = -2 \sum_{i=1}^n y_i \log(\hat{p}_i) + (1 - y_i) \log(1 - \hat{p}_i)$$

- We're left with the loglikelihood summed up over all the points, times -2. If you don't calculate the deviance itself, you just look at the loglikelihood for a given point and that's that some value, that's the residual.

- The deviance and the loglikelihood are equal to each other, they are the same. This is also the same formula for the logarithmic scoring.
- When we talk about accuracy, the likelihood is a measure of accuracy. Sometimes we transform it in different ways and sometimes we leave it the same way and call it by a different name, because it's just very good a measure of accuracy.

Deviance

```
summary(lmod)
```

```
##
## Call:
## glm(formula = chd ~ height + cigs, family = binomial, data = wcgs)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.0041  -0.4425  -0.3630  -0.3499   2.4357
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -4.50161    1.84186  -2.444  0.0145 *
## height      0.02521    0.02633   0.957  0.3383
## cigs        0.02313    0.00404   5.724 1.04e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 1781.2  on 3153  degrees of freedom
## Residual deviance: 1749.0  on 3151  degrees of freedom
## AIC: 1755
##
## Number of Fisher Scoring iterations: 5
```

- The null deviance and the residual deviance are given, with large degrees of freedom.
- The residual deviance is the deviance for this value, it has 3151 degrees of freedom because it's the number observations minus the number of parameters (three, including the intercept)
- The null deviance is if you don't have anything in the model, so you don't have the height and cigs and just have the intercept
- The difference between two deviances follow the chi squared distribution. So we can set up a hypothesis test. We can say that our two models with different parameters where one is the subset of the other, are they different in how they work? Same thing as saying should we keep height in the model?

Deviance

```
p_hat <- fitted.values(lmod)

dev_res <- -2*sum(p_hat*logit(p_hat) + log(1-p_hat))

null_mod <- update(lmod, formula. = . ~ 1)
p_hat_null <- fitted.values(null_mod)
dev_null <- -2*sum(p_hat_null*logit(p_hat_null) + log(1-p_hat_null))

dev_res
```

```
## [1] 1749.049
```

```
dev_null
```

```
## [1] 1781.244
```

Compare two models

- Degrees of freedom is the difference in df between the two models
- You can do a chi square test between those two things and say if those two models are different or not
- The 2 is the difference in degrees of freedom of 3153-3151, which is 2.

```
1 - pchisq(dev_null - dev_res, 2)
```

```
## [1] 1.02106e-07
```

- We get a very low p-values, this suggests that the 2 variables we included are ‘useful’ then the model with just the intercept

Look at a single variable

- You can do the same thing using the anova function and it does the same exact thing.
- Now we’re going to test for height. We’re going to create another model which excludes height, and just has cigs.

```
lmod <- glm(chd ~ height + cigs, family = binomial, wcsv)
lmodc <- glm(chd ~ cigs, family = binomial, wcsv)
anova(lmodc, lmod, test="Chi")
```

Resid. Df	Resid. Dev	Df	Deviance	Pr(>Chi)
3152	1749.969	NA	NA	NA
3151	1749.049	1	0.9202473	0.3374101

- The p value is greater than .05, suggesting that we don’t need height in the model. We can drop out height from the model
- You always want to keep in mind what do you want out of this model
- If it’s inference sometimes. you want to keep stuff in. If it’s prediction maybe you want to drop stuff out.
- People generally like to have similar models because the idea is that you’re less likely to over figure data, and it’s better to generalize better
- If you really want to keep something in just because it’s not significant to the model doesn’t mean that its bad to have it there
- That is how to test the inclusion of a variable ^

Confidence Interval (CI)

$$\hat{\beta}_i \pm z^{\alpha/2} se(\hat{\beta}_i)$$

For β_1 (height)

- You can test the confidence intervals using the normal distribution
- The estimate + alpha .025 is the 1.96* by the std error

```
exp(0.02521 + c(-1,1) * 1.96 * 0.02633)
```

```
## [1] 0.9739486 1.0798442
```

```
confint(lmod)
```

```
## Waiting for profiling to be done...
```

	2.5 %	97.5 %
(Intercept)	-8.1347547	-0.9129702
height	-0.0261990	0.0770283
cigs	0.0151495	0.0310053

- The CI for the beta is between .97 and 1.07 is exponentiated so it's an odds ratio.
- Really what we care about odds ratio is if it crosses 1 or not. If it's greater than 1, that means it has a positive outcome. If it's less than 1, it has less of an effect.
- Think you have to load the mass library first, but it seems to work
- It looks at the likelihood function and estimates the CI from that, here they are still on a linear scale not exponential.
- These are going to be a little wider than the first equation. For sake of argument, $\exp(-0.026) = 0.974$. While the first one was .973
- Likelihood is generally better, its more accurate

Residuals

- The way you can calculate residuals is that if it's a yes it's a 1, and you subtract the predicted probability that's a residual.

```
pred_prob <- predict(lmod, type="response")
raw_res <- as.numeric((wcgs$chd == "yes")) - pred_prob
raw_res <- residuals(lmod, type = "response")
```

- The problem though is that we can't do same thing with residuals that we can do with a linear model because the variance with this residual is not constant. It's higher when you're closer to .5 and lower when you're at 1 or 0. Because of this, it is hard to gauge homogeneity

Residuals' Variance

- Variance is not constant (binary variance = $p(1-p)$)
- What we can do is we can actually use the deviance, which is the square root value in the equation below and make it just negative or positive
- Use deviance residual instead

$$r_i = \text{sign}(y_i - \hat{p}_i) \sqrt{d_i^2}$$

- The deviance residuals are not constrained to have mean zero so the mean level in the plot is not of interest.

```
wcgs2 <- wcgs2 %>%
  mutate(residuals=residuals(lmod), linpred=predict(lmod))
```

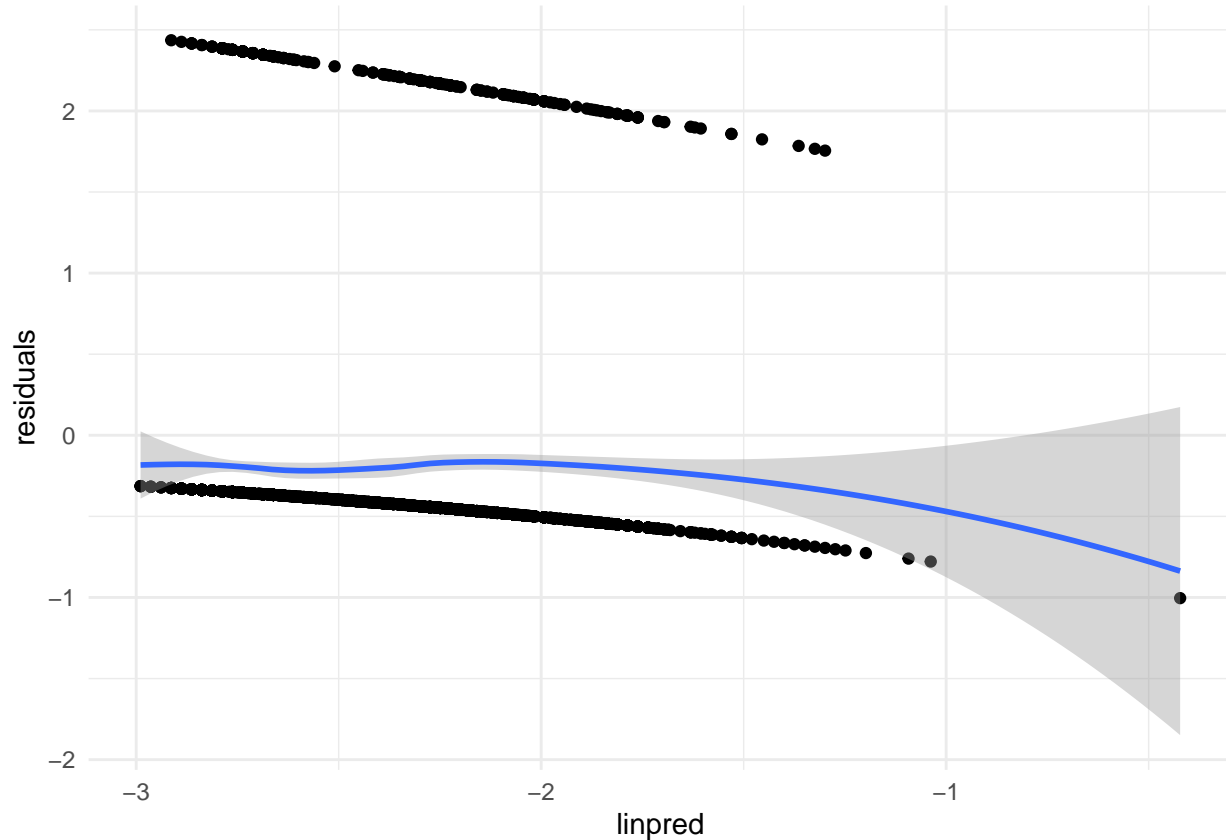
Plots of residuals

- Use smoothing instead of binning that's in the book
- Graph the linear predictor (eta) on the x-axis against the residuals on the y-axis.
- The textbook likes to bin them together, say cut that into 10 pieces and calculate the average and look at it that way. I don't think that's a good idea

- Instead put a smoothing function on it. What you're looking for is some sort of deviation from a pattern.

```
wcgs2 %>%
  ggplot(aes(x = linpred, y = residuals)) +
  geom_point() +
  geom_smooth(method = 'loess')
```

```
## `geom_smooth()` using formula 'y ~ x'
```



- What we see here is that it's pretty straight. Unlike previously, the residuals don't have to add up to 0 so it's okay that it's sort of below 0.

Plots of residuals

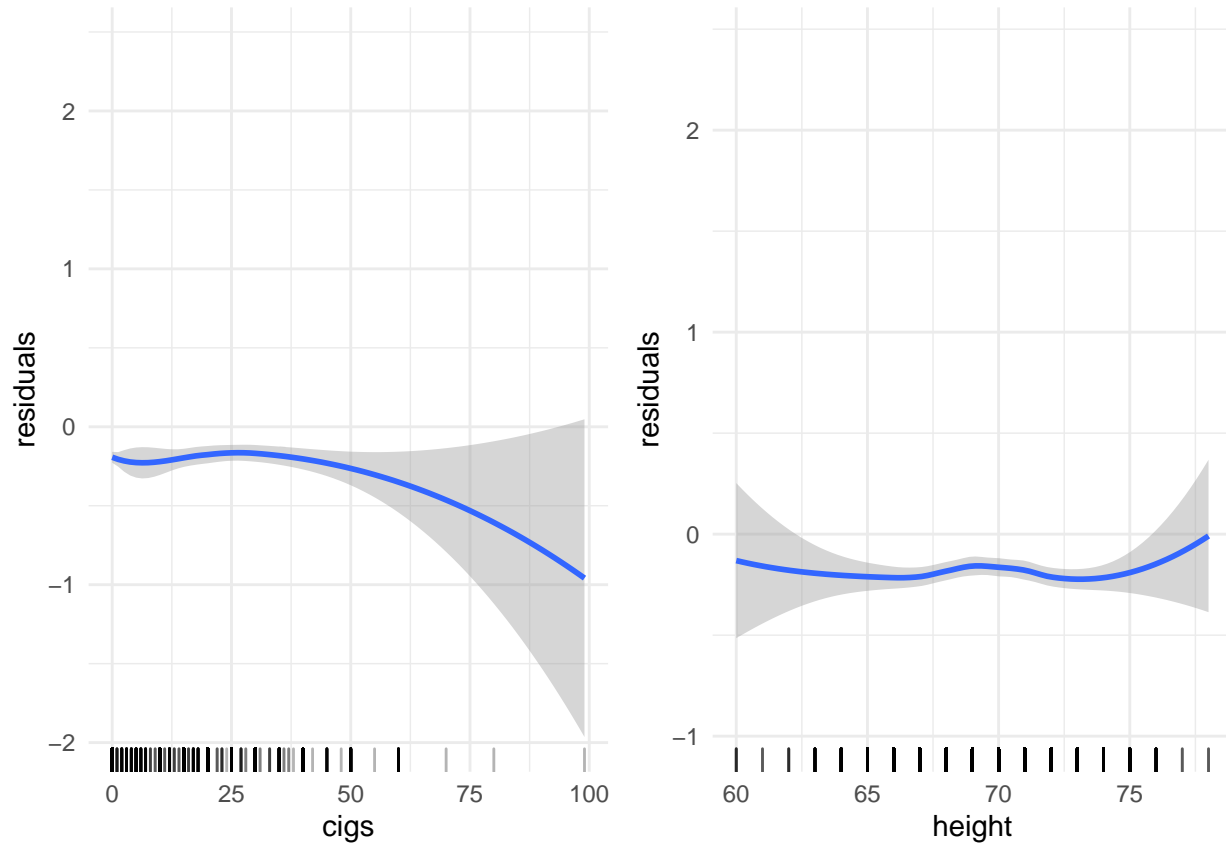
- You can do the same thing but use individual predictors

```
p1 <- wcgs2 %>%
  ggplot(aes(x = cigs, y = residuals)) +
  # geom_point() +
  geom_smooth(method = 'loess') +
  geom_rug(sides = "b", alpha = .3)

p2 <- wcgs2 %>%
  ggplot(aes(x = height, y = residuals)) +
  # geom_point() +
  geom_smooth(method = 'loess') +
  geom_rug(sides = "b", alpha = .3)
```

```
grid.arrange(p1, p2, ncol = 2)
```

```
## `geom_smooth()` using formula 'y ~ x'
## `geom_smooth()` using formula 'y ~ x'
```



- This looks okay, if you saw any crazy patterns doing up or down. We see this it's pretty flat and straight

Leverage

```
wcgs2 %>%
  mutate(lev = hatvalues(lmod)) %>%
  arrange(-lev) %>%
  dplyr::select(height, cigs, lev)
```

height	cigs	lev
71	99	0.0280034
64	80	0.0187294
76	60	0.0107423
71	70	0.0098510
74	60	0.0080759
66	60	0.0072673
73	60	0.0071473
67	60	0.0066109
72	60	0.0064751
72	60	0.0064751
74	55	0.0064741

height	cigs	lev
60	30	0.0063419
75	50	0.0061798
64	50	0.0061661
68	60	0.0061509
68	60	0.0061509
71	60	0.0060494
71	60	0.0060494
62	40	0.0060221
69	60	0.0058973
70	60	0.0058602
70	60	0.0058602
74	50	0.0051461
66	50	0.0045425
73	50	0.0043484
73	50	0.0043484
73	50	0.0043484
73	50	0.0043484
73	50	0.0043484
65	45	0.0042372
64	40	0.0041635
75	40	0.0040919
75	40	0.0040919
75	40	0.0040919
60	10	0.0040861
74	45	0.0040676
67	50	0.0039639
76	31	0.0038413
72	50	0.0037768
72	50	0.0037768
72	50	0.0037768
72	50	0.0037768
72	50	0.0037768
72	50	0.0037768
76	30	0.0037257
78	0	0.0036913
78	0	0.0036913
78	0	0.0036913
68	50	0.0035529
68	50	0.0035529
68	50	0.0035529
68	50	0.0035529
68	50	0.0035529
60	0	0.0034690
60	0	0.0034690
60	0	0.0034690
64	35	0.0034520
71	50	0.0034213
71	50	0.0034213
65	40	0.0034051
72	48	0.0033617
73	45	0.0033328
73	45	0.0033328

height	cigs	lev
69	50	0.0033191
69	50	0.0033191
69	50	0.0033191
70	50	0.0032720
70	50	0.0032720
70	50	0.0032720
70	50	0.0032720
70	50	0.0032720
70	50	0.0032720
74	40	0.0032118
74	40	0.0032118
74	40	0.0032118
74	40	0.0032118
74	40	0.0032118
74	40	0.0032118
74	40	0.0032118
74	40	0.0032118
77	5	0.0030354
64	31	0.0029969
61	0	0.0029424
61	0	0.0029424
61	0	0.0029424
77	0	0.0029086
77	0	0.0029086
64	30	0.0028973
64	30	0.0028973
76	20	0.0028681
76	20	0.0028681
76	20	0.0028681
62	10	0.0028239
75	30	0.0027992
75	30	0.0027992
75	30	0.0027992
66	40	0.0027712
66	40	0.0027712
66	40	0.0027712
66	40	0.0027712
66	40	0.0027712
66	40	0.0027712
63	20	0.0027665
63	20	0.0027665
63	20	0.0027665
63	20	0.0027665
65	35	0.0027548
65	35	0.0027548
63	18	0.0026428
74	35	0.0025509
74	35	0.0025509
73	40	0.0025369
73	40	0.0025369
73	40	0.0025369
73	40	0.0025369
73	40	0.0025369

[illegible]

[illegible]

[illegible]

[illegible]

height	cigs	lev
69	10	0.0003717
69	10	0.0003717
69	10	0.0003717
69	10	0.0003717
69	10	0.0003717
69	10	0.0003717
69	10	0.0003717
69	10	0.0003717
69	10	0.0003717
69	10	0.0003717
69	10	0.0003717
69	10	0.0003717
69	10	0.0003717
70	7	0.0003697
70	7	0.0003697
70	7	0.0003697
70	7	0.0003697
70	7	0.0003697
69	15	0.0003679
69	15	0.0003679
69	15	0.0003679
69	15	0.0003679
69	15	0.0003679
69	15	0.0003679
69	15	0.0003679
69	15	0.0003679
69	15	0.0003679
69	15	0.0003679
69	15	0.0003679
69	15	0.0003679
69	12	0.0003638
69	12	0.0003638
69	12	0.0003638
69	12	0.0003638
70	18	0.0003544
70	9	0.0003488
70	17	0.0003439
70	17	0.0003439
70	17	0.0003439
70	10	0.0003408
70	10	0.0003408
70	10	0.0003408
70	10	0.0003408
70	10	0.0003408
70	10	0.0003408
70	10	0.0003408
70	10	0.0003408
70	10	0.0003408
70	10	0.0003408
70	10	0.0003408
70	10	0.0003408
70	16	0.0003362

height	cigs	lev
70	15	0.0003312
70	15	0.0003312
70	15	0.0003312
70	15	0.0003312
70	15	0.0003312
70	15	0.0003312
70	15	0.0003312
70	15	0.0003312
70	15	0.0003312
70	15	0.0003312
70	15	0.0003312
70	15	0.0003312
70	15	0.0003312
70	15	0.0003312
70	15	0.0003312
70	15	0.0003312
70	12	0.0003305
70	12	0.0003305
70	12	0.0003305
70	14	0.0003286
70	13	0.0003285
70	13	0.0003285
70	13	0.0003285

Model Selection

- Using chi square for model selection is not super great. it just doesn't really work how it should
- One way we could do it that's better is to look at the AIC
- $AIC = (-2\log\text{likelihood} + 2q)$, 2 is the number of parameters
- This is a way to penalize adding extra parameters to the data
- That means it could be a better way of judging whether something is fitting well

$$AIC = -2\log L + 2q$$

- Create a huge model, with many parameters.
- The step function in R looks at each one of these and removes a variable one at a time and refits the model
- If you remember weight, here's the AIC and etc. It's going to keep doing this until the none variable is at the top and then that is your final model
- Meaning that it's not dropping another after that has the lowest AIC

```
wcgs2$bmi <- with(wcgs2, 703*weight/(height^2))
lmod <- glm(chd ~ age + height + weight + bmi + sdp + dbp + chol + dibep + cigs + arcus, family=binomial,
lmodr <- step(lmod, trace=1)
```

```
## Start:  AIC=1591.21
## chd ~ age + height + weight + bmi + sdp + dbp + chol + dibep +
##      cigs + arcus
##
##           Df Deviance    AIC
## - dbp      1   1569.2 1589.2
## - weight   1   1569.3 1589.3
## - bmi      1   1569.5 1589.5
```

```
## - height 1 1569.5 1589.5
## <none> 1569.2 1591.2
## - arcus 1 1571.3 1591.3
## - sdp 1 1577.0 1597.0
## - dibep 1 1590.5 1610.5
## - cigs 1 1592.2 1612.2
## - age 1 1593.8 1613.8
## - chol 1 1620.0 1640.0
##
## Step: AIC=1589.22
## chd ~ age + height + weight + bmi + sdp + chol + dibep + cigs +
## arcus
##
## Df Deviance AIC
## - weight 1 1569.3 1587.3
## - bmi 1 1569.5 1587.5
## - height 1 1569.5 1587.5
## <none> 1569.2 1589.2
## - arcus 1 1571.3 1589.3
## - sdp 1 1586.5 1604.5
## - dibep 1 1590.5 1608.5
## - cigs 1 1592.7 1610.7
## - age 1 1593.8 1611.8
## - chol 1 1620.0 1638.0
##
## Step: AIC=1587.33
## chd ~ age + height + bmi + sdp + chol + dibep + cigs + arcus
##
## Df Deviance AIC
## <none> 1569.3 1587.3
## - arcus 1 1571.5 1587.5
## - height 1 1572.6 1588.6
## - bmi 1 1574.4 1590.4
## - sdp 1 1586.6 1602.6
## - dibep 1 1590.7 1606.7
## - cigs 1 1592.8 1608.8
## - age 1 1593.9 1609.9
## - chol 1 1620.0 1636.0
```

- In this case, it drops dbp and weight. What that says is that these 2 variables are not important variables to the model given everything else that is in it.

Resulting Model

- It keeps some stuff that's not statistically significantly different from 0, looking at the p-value. The arcus present is not significant but nevertheless it still kept it in the model using this method
- Here you don't need to drop out all the variable to "create the best model"
- This is ok if your goal is future *prediction* of data

```
summary(lmodr)
```

```
##
## Call:
## glm(formula = chd ~ age + height + bmi + sdp + chol + dibep +
## cigs + arcus, family = binomial, data = wcgs2)
```

```
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.3225  -0.4352  -0.3118  -0.2199   2.8517
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -15.957599   2.286076  -6.980 2.94e-12 ***
## age          0.061590   0.012397   4.968 6.76e-07 ***
## height       0.050161   0.027824   1.803  0.0714 .
## bmi          0.060385   0.026599   2.270  0.0232 *
## sdg          0.017728   0.004155   4.267 1.98e-05 ***
## chol         0.010709   0.001529   7.006 2.45e-12 ***
## dibepB       0.657616   0.145898   4.507 6.56e-06 ***
## cigs         0.021041   0.004262   4.936 7.96e-07 ***
## arcuspresent  0.210998   0.143718   1.468  0.1421
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 1769.2  on 3139  degrees of freedom
## Residual deviance: 1569.3  on 3131  degrees of freedom
## (14 observations deleted due to missingness)
## AIC: 1587.3
##
## Number of Fisher Scoring iterations: 6
```

But what if you want inference?

- However you can go ahead and say that dbp has nothing to do with chd just because it was left out of the model. Because if you look at it individually, you'll see that there is a very strong significance. What this probability means is that its correlated with one of these other variables and you basically don't need both of them in the model.
- Just cause a factor is dropped from the model doesn't mean it may not influence outcome
- dbp was dropped from the model but that doesn't mean it doesn't influence the outcome
- You always have to balance what you're trying to do in terms of inference and prediction of data

```
summary(update(lmod, . ~ dbp))
```

```
##
## Call:
## glm(formula = chd ~ dbp, family = binomial, data = wgs2)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.0988  -0.4298  -0.3902  -0.3541   2.4976
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -5.221713   0.511646 -10.206 < 2e-16 ***
## dbp          0.033560   0.005981   5.611 2.01e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
```

```
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 1781.2 on 3153 degrees of freedom
## Residual deviance: 1751.7 on 3152 degrees of freedom
## AIC: 1755.7
##
## Number of Fisher Scoring iterations: 5
```

- The stars and dots in the significant part: the *** mean that the p-value is effectively 0, the . means that it's between .05 and 1 and nothing means that its above 0.1.

Goodness of Fit

Does a model fit?

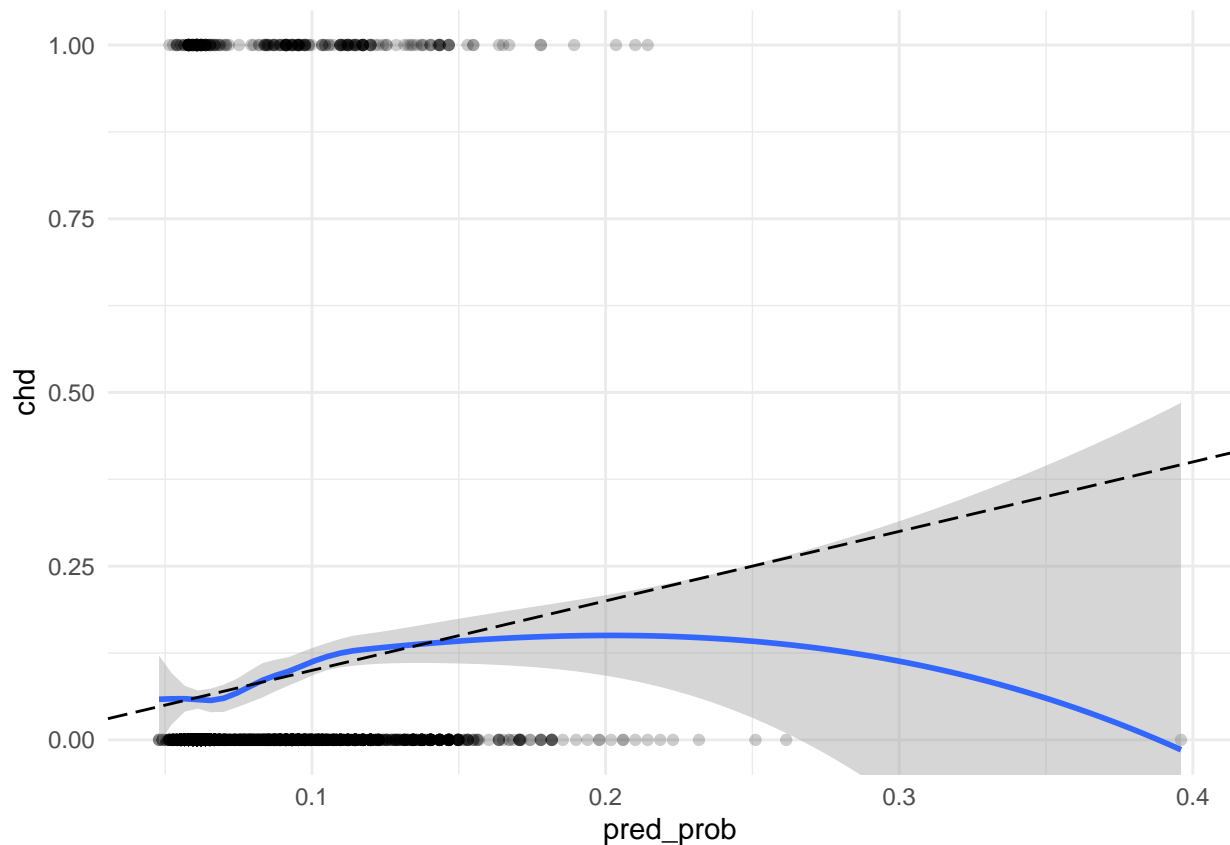
- I don't think this is a particularly useful question
- I've never come across a problem where a model doesn't fit at all

Calibration Curve

- The single most useful graph for a binary variables is a calibration curve
- This curve plots on the x-axis the predictions and the y-axis it plots the response or true value there.
- The dots is not very interesting so you apply a smoothing function to it, and what you want is for your model to be well calibrated. Calibrated means that every time I assign a probability of 10% of chd, if you look at 100 people who I have assigned 10% to ideally 10% of them would get the disease meaning 10 people would have it.
- That's what this curve is showing. The dotted line goes through 0,0 and has a slope of 1. So you want the curve to basically fall along that dotted line.
- It starts to dip down towards the right, but you'll also notice that we have a lot less data there and that the CI of the smoothing curve is very wide
- Being below this line suggests that we are over estimating the risk of chd in the population
- Example: look at the .2 pred_prob which says that this person has a 20% chance of having chd, but in reality maybe only 17% do
- This is a good way of diagnosing and looking at binary fits of data and convincing myself that the model is good
- We would always want it to follow that diagonal linear dotted line
- The textbook shows a binned version (figure 2.9), where they take everybody who has a probability of 5-15% and looks at the average of occurrences of chd in that group. This method has lots of times very weird stuff happen for no reason other than the fact that maybe you didn't choose good break points.

```
wcgs2 %>%
  ggplot(aes(x = pred_prob, y = chd)) +
  geom_point(alpha = .2) +
  geom_smooth(method = "loess") +
  geom_abline(slope = 1, intercept = 0, linetype = "longdash") +
  coord_cartesian(ylim = c(0,1))
```

```
## `geom_smooth()` using formula 'y ~ x'
```

Scoring

- Another way to think about calibration curve
- Two type of scores you can give an algorithm: brier and logarithmic
- The idea is to summarize the performing algorithm in one number
- Brier Score (nicer to interpret)
- The average of the difference between your prediction and the outcome, squared.
- It's a nice way to look at it
- For example, professor used this score to placed bets or predictions on last avengers movie on who would be alive at the end of the movie.

$$\frac{1}{n} \sum_{i=1}^n (\hat{p}_i - y_i)^2$$

- logarithmic scoring (okay)

$$\sum y \log(\hat{p}) + (1 - y) \log(1 - \hat{p})$$

- In both of these types of scores, you get a higher score the more correct you are. So if I say something is 90% likely to happen and then it happens, I should be considered “more right” than somebody who said that there’s a 50% chance that it’ll happen and it happens.
- In both of these cases you get that play. But the brier score is a little bit more interpretable than the logarithmic, and it’s used a lot in election forecasts, etc.

Confusion Matrix

```
tabl <- wgs %>%  
  filter(complete.cases(.)) %>%  
  mutate(predout = ifelse(fitted.values(lmodr) > .5, "yes", "no")) %>%  
  xtabs(~chd + predout, .)
```

tabl

chd/predout	no	yes
no	2882	3
yes	253	2

- This is called a confusion matrix, or a two by two square. It shows you the different combinations of the actual label (two rows on left side) and the predicted (from the model) label (two columns on right side). The predictions are no and yes, and actual values are no and yes.
- When we're taking about how well a model is performing, one of the easiest ways to think about it often is just look at the accuracy of the model.
- To find accuracy, what we do is we add up the number of correct cases and divide it by the total number of cases.
- Confusion matrix is good to have to see the yes or no labels to understand the numbers
- Accuracy =

```
(2882+2)/(2882+3+253+2)
```

```
## [1] 0.9184713
```

- 91% accurate.
- Accuracy doesn't tell the full story. One of the things that you notice immediately is the 253. We only predicted 5 people to have chd, but we know that's not really accurate in terms of the population. What we're seeing in the 253 is that we're missing a lot of ppl, saying no you don't have chd, but in fact you do.
- Accuracy can be misleading, especially in an imbalanced dataset
- An imbalanced dataset is one where the number of people where you have the thing you're trying to predict is not very high. In this case, it's about 10%

Specificity and Sensitivity

Specificity (for those that don't have the disease how many were predicted to not have it) (chd = no group)

tabl

chd/predout	no	yes
no	2882	3
yes	253	2

- That is of the people who do not have the disease, which is the no row, how many did we predict correctly, which is the 2882 value. To find that, you do the equation below

```
2882/(2882 + 3)
```

```
## [1] 0.9989601
```

- This number is very high, 99% is the specificity. If we get every single person who is a no correct, we'll

have 100% specificity.

Sensitivity (for those that have disease how many were correctly identified?) (chd = yes) - That number is shown below

```
2/(253 + 2)
```

```
## [1] 0.007843137
```

- This number is very low, that goes back to this idea that we weren't actually identifying the chd candidates very well because of all the people that were in 253.
- That is a way of summarizing the problem we have with the accuracy test
- Specificity and sensitivity only pertains when you have a binary value. It really applies to any sort of time that you're trying to estimate another quantity. Like if you're trying to build and advertising algorithm, how well do you predict that this person is going to buy the product. You'll be able to put that in to a confusion matrix as well and you'll be able to calculate the same things there
- The sensitivity is the positive case, like if you're sensitive to something, and the specificity is the opposite of that
- These two specificity and sensitivity are weird values, because you're conditioning on something that you don't know. You don't know if somebody doesn't have chd when you're putting this test together, you only find out later through some other things. So you're sort of calculating the statistic without knowing so well who these ppl are

PPV and NPV

- One of the other ways of looking at it is the positive predicted value and the negative predicted value.

```
tabl
```

chd/predout	no	yes
no	2882	3
yes	253	2

- It may be more reasonable to ask: What's the probability that I have the disease if I test positive?

PPV - The PPV is more in line with if we're talking about this algorithm or a test of some sort, and I get a yes, how likely am I to have chd. You would do that by looking at everyone you predicted yes for, and how accurate you were when you predicted yes

```
2/(3 + 2)
```

```
## [1] 0.4
```

NPV - The NPV is when you say no, how many of them actually end up being no - Same thing as, the probability if you test negative, you don't have the disease

```
2882/(253 + 2882)
```

```
## [1] 0.9192982
```

- It's a different way to look at it. People always talk about specificity and sensitivity, however when you actually present the results of a model, whether they know it or not, I think what they care about is PPV and NPV.
- What we sort of implying?

```
tabl <- wgs %>%  
  filter(complete.cases(.)) %>%
```

```
mutate(predout = ifelse(fitted.values(lmodr) > .5, "yes", "no")) %>%
xtabs(~chd + predout, .)
```

tabl

chd/predout	no	yes
no	2882	3
yes	253	2

- Our logistic regressions put out a probability value, a number between 0 and 1. The way I created these labels of knowing yes, I had to pick a cut off value. So how do you do that, you have to choose a number, in this case what's often used by default is 0.5.
- But we can actual vary that. What if we say we actually want to label anyone with 20% or 10% higher as a yes instead of using the 0.5 cut off.

Vary the cutoff

- We can do that by:
- What's going on here is that we're looking at a few thresholds, we're sequencing from 0.01 to 0.5, at a 0.01 interval.
- And I'm creating a bunch of these things, sensitivity, specificity, PPV, NPV and F1 score. I'm basically creating those tables and calculating all of those values here in a for-loop.
- And I just put all of it in one big data frame.
- What you can see is there are relationships between these values

```
wcgs2 <- wcgs2 %>%
  filter(complete.cases(.)) %>%
  mutate(predprob = predict(lmodr, ., type = "response"))

thresh <- seq(0.01, 0.5, 0.01)
Sensitivity <- numeric(length(thresh))
Specificity <- numeric(length(thresh))
PPV <- numeric(length(thresh))
NPV <- numeric(length(thresh))
F1 <- numeric(length(thresh))
for(j in seq(along=thresh)){
  pp <- ifelse(wcgs2$predprob < thresh[j], "no", "yes")
  xx <- xtabs(~ chd + pp, wcgs2)
  Specificity[j] <- xx[1,1]/(xx[1,1]+xx[1,2])
  Sensitivity[j] <- xx[2,2]/(xx[2,1]+xx[2,2])
  PPV[j] <- xx[2,2]/(xx[1,2] + xx[2,2])
  NPV[j] <- xx[1,1]/(xx[2,1] + xx[1, 1])
  F1[j] <- xx[2,2]/(xx[2,2] + 1/2*(xx[2,1] + xx[1,2]))
}
full_df <- data.frame(
  thresh = thresh,
  sens = Sensitivity,
  spec = Specificity,
  PPV = PPV,
  NPV = NPV,
  F1 = F1
)
```

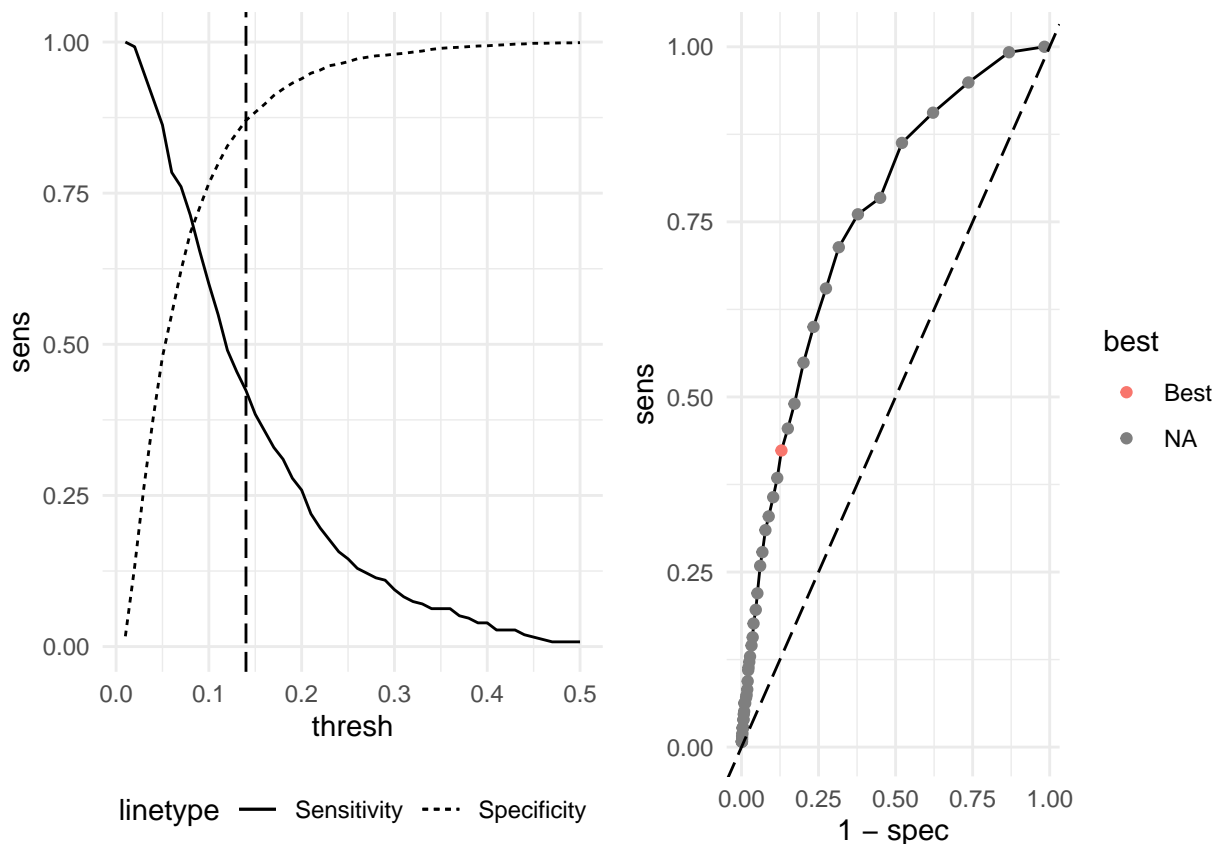
```
full_df$best <- NA
full_df$best[which.max(F1)] <- "Best"
best_thresh <- thresh[which.max(F1)]
```

Sensitivity vs Specificity

```
p1 <- ggplot(full_df, aes(x = thresh)) +
  geom_line(aes(y = sens, linetype = "Sensitivity")) +
  geom_line(aes(y = spec, linetype = "Specificity")) +
  theme(legend.position = 'bottom') +
  geom_vline(xintercept = best_thresh, linetype = "longdash")

p2 <- ggplot(full_df, aes(x = 1-spec, y = sens)) +
  geom_line() +
  geom_point(aes(color = best)) +
  geom_abline(slope = 1, intercept = 0, linetype = "longdash")

grid.arrange(p1, p2, ncol = 2)
```



Don't pay attention to the labels y-axis, look at the line type below.

Left graph

- Here we're looking at specificity and sensitivity. So if you look at the threshold here, 0.1-0.5. what is the specificity and sensitivity of using a particular threshold?
- What happens is when you have a high sensitivity, you have a low specificity. and as you move along

the threshold, those two things start to switch spots. And the trick is somehow choosing a point where they are a reasonable trade off between the two values.

- When we add a threshold value, we're basically adding another parameter to the model. Should we validate our threshold choice against the held out set? We're not going to cover that in this class because there's a whole other set of theories of cross-validating and other stuff. But in the same way you can overfit your data by including too many parameters, you could also overfit your data by choosing a threshold. Here though we're not going to worry about that now. For now, we're just going to pick a threshold based off what we see in this data.
- The intersection between the two is unfortunately rarely used as any point because it doesn't really mean anything.
- The dotted line straight down is the threshold that is "optimal" by one specific characteristic.
- When it sums down to it, it's sort of a business decision. You have to help people understand what it means, do you want more false positives or do you want to miss people? That's what the tradeoff is going to be.
- In this case, chd, what if this is just some screening algorithm that just runs in the background. Every time a positive value gets hit for chd, you have to call the person get them in for an appointment, and run a battery of tests on them. That can be a very expensive thing that takes a lot of medical staff time and the patient's time, and you could be finding too many people that way. And if you find this person with chd, maybe you find them 3 months earlier than otherwise because maybe their coming in for an appointment or something. In which case you could say I don't really want any false positives, I'm cool with missing a lot of ppl, but I want to make sure that whenever I do get a correct prediction I'm really sure about it. In that case, you would want a really high specificity. Looking at the matrix, I want to make sure that I don't have a lot of people who I say are predicted as having chd but end up not having it. You want the no row and the yes column number to be low. And if that is low, that means the no row and no column number should be little higher. And that is saying we want a high specificity there.
- A low value is not necessarily a bad thing, it is based off all these sort of decisions and it's okay for an algorithm to have a very low threshold.
- You can guess the threshold, you can go from 0.1 to 1, and you can just look at it. The reason why we only go to .5 in this case if you look at the calibration curve, is that we only predict one 1 person above .5, so there's no point at looking higher.

Right graph

- This is your classic, receiver operating characteristic (ROC) curve. This is often reported alongside binary classification models.
- What this is that it's not showing the threshold anymore. What it's showing is 1- specificity on the x-axis and the sensitivity on the y-axis.
- Each one of these dots on the curve part is a different threshold, which lends itself to a different combination of specificity and sensitivity.
- If your model does absolutely nothing, basically just guesses based on the rate of the population, your ROC will follow the dashed straight line here, it'll be one at a 45 degree angle. If it does worse, it'll run on the bottom side of that straight line. That would only happen if you overfit that data.
- What you want is a really high specificity and a really high sensitivity, that would be a perfect model. The way you could do that is you could have a point at the 1 of the y-axis, if you have a point there that would mean that you have a perfect model. You sort want to be up here as much as possible, it rarely happens, but that is what you would want to see.
- This is also the curve, the area under the curve of the ROC (AUC). The area under the ROC curve. A perfect AUC is 1, and 50% is random. The higher your AUC, the better your algorithm is performing.
- The F1 score is the balance of specificity and sensitivity. You can calculate it and the highest F1 is maybe your best threshold.

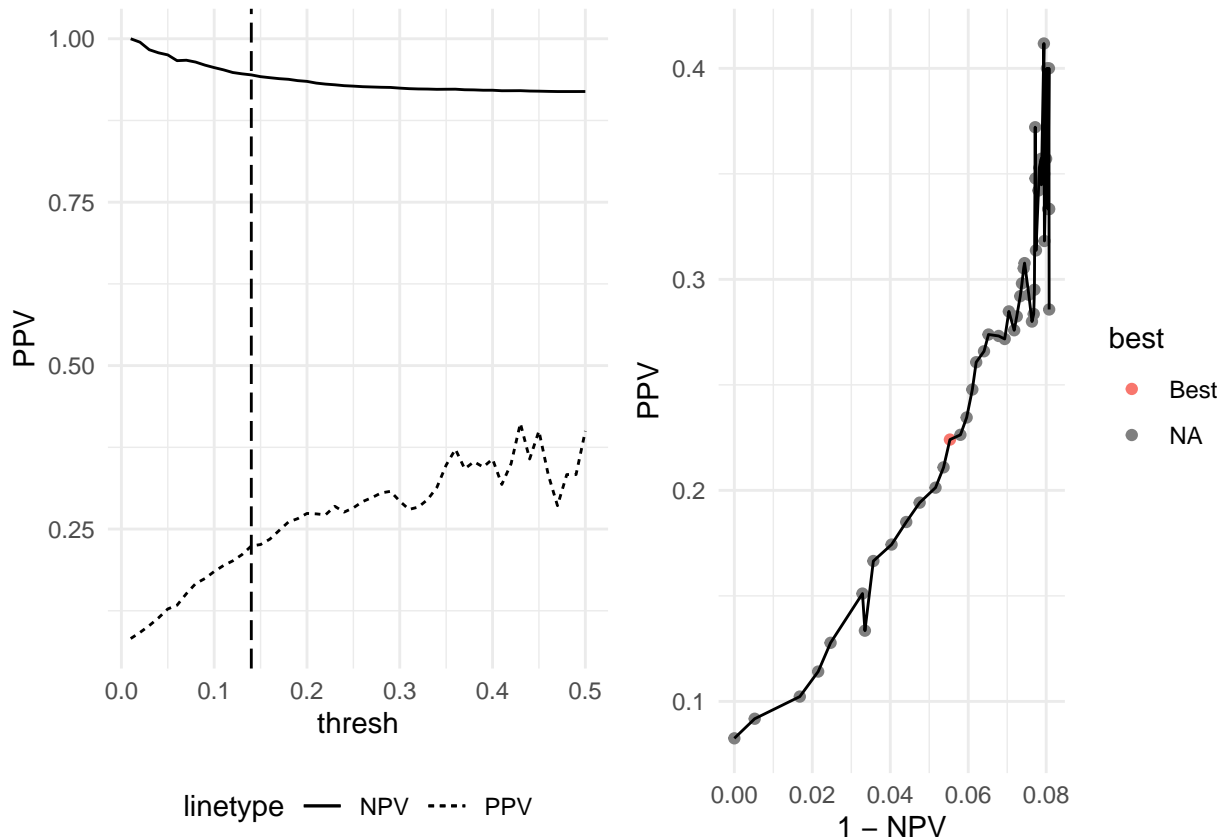
- In these two graphs, the F1 score shows that in the right, the red dot is the best " on ". And the left, it's the straight dotting line, which is .14.

NPV and PPV

```
p1 <- ggplot(full_df, aes(x = thresh)) +
  geom_line(aes(y = PPV, linetype = "PPV")) +
  geom_line(aes(y = NPV, linetype = "NPV")) +
  theme(legend.position = 'bottom') +
  geom_vline(xintercept = best_thresh, linetype = "longdash")

p2 <- ggplot(full_df, aes(x = 1 - NPV, y = PPV)) +
  geom_point(aes(color = best)) +
  geom_line()

grid.arrange(p1, p2, ncol = 2)
```



- Again, don't pay attention to the labels y-axis, look at the line type below.
- Something that is tough with the NPV and the PPV, is the curves aren't quite as easy to look at and interpret.
- The NPV goes down and the PPV goes up as you change the threshold.
- It is just the graph but it hard to see what it's going on.
- But both show the best value is the red dot on the right, and the on the threshold as well the straight dotted line on the left.

Not covered

- R^2 - in terms of the binomial, it's sort of the adaption of the R^2 in the linear model. Not interesting and not very useful. With binomial data, your value of R^2 shouldn't really be the end all of in your data analysis.
- So because of that, we should look at the specificity and sensitivity curve, and the calibration curve, the best ways to look at it and see if we're heading in the right direction. Then you have to think about it if you care about accuracy or something else, it's just a little bit harder in the Bernoulli distribution on the binary outcome. It's hard because you don't have this thing to look at, a nice distribution of varying numbers to look at it, so it can be hard to summarize things.
- Something that can help is the scores.
- Estimation issues - if you have perfect separation between a variable and an outcome, that'll cause issues with the actual algorithm that's fitting it. Read about it in textbook, but it's never happened in professor's life.

Lab with pima dataset (question 3 of exercise)

```
summary(pima)
```

pregnant	glucose	diastolic	triceps	insulin	bmi	diabetes
Min. : 0.000	Min. : 0.0	Min. : 0.00	Min. : 0.00	Min. : 0.0	Min. : 0.00	Min. :0.0780
1st Qu.: 1.000	1st Qu.: 99.0	1st Qu.: 62.00	1st Qu.: 0.00	1st Qu.: 0.0	1st Qu.:27.30	1st Qu.:0.2437
Median : 3.000	Median :117.0	Median : 72.00	Median :23.00	Median : 30.5	Median :32.00	Median :0.3725
Mean : 3.845	Mean :120.9	Mean : 69.11	Mean :20.54	Mean : 79.8	Mean :31.99	Mean :0.4719
3rd Qu.: 6.000	3rd Qu.:140.2	3rd Qu.: 80.00	3rd Qu.:32.00	3rd Qu.:127.2	3rd Qu.:36.60	3rd Qu.:0.6262
Max. :17.000	Max. :199.0	Max. :122.00	Max. :99.00	Max. :846.0	Max. :67.10	Max. :2.4200

```
mod <- glm(cbind(test, 1 - test) ~ ., data = pima, family = "binomial")
summary(mod)
```

```
##
## Call:
## glm(formula = cbind(test, 1 - test) ~ ., family = "binomial",
##      data = pima)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.5566  -0.7274  -0.4159   0.7267   2.9297
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -8.4046964  0.7166359 -11.728  < 2e-16 ***
## pregnant     0.1231823  0.0320776   3.840 0.000123 ***
## glucose      0.0351637  0.0037087   9.481  < 2e-16 ***
## diastolic    -0.0132955  0.0052336  -2.540 0.011072 *
## triceps      0.0006190  0.0068994   0.090 0.928515
## insulin     -0.0011917  0.0009012  -1.322 0.186065
## bmi          0.0897010  0.0150876   5.945 2.76e-09 ***
## diabetes     0.9451797  0.2991475   3.160 0.001580 **
## age          0.0148690  0.0093348   1.593 0.111192
## ---
```



```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##    Null deviance: 993.48  on 767  degrees of freedom
## Residual deviance: 723.45  on 759  degrees of freedom
## AIC: 741.45
##
## Number of Fisher Scoring iterations: 5
pchisq(723.45, 759, lower = FALSE)

## [1] 0.8185642
pima %>% filter(glucose == 99, pregnant ==1 )
```

pregnant	glucose	diastolic	triceps	insulin	bmi	diabetes	age	test
1	99	72	30	18	38.6	0.412	21	0
1	99	58	10	0	25.4	0.551	21	0