

# Predictive Analytics Lecture 3

---

Adam Kapelner

Stat 422/722

at The Wharton School of the University of Pennsylvania

---

January 31 & February 1, 2017

## The Coin Example from Last Class I

I want to explain the coin example from last class in the context of likelihood. Imagine you flip a coin three times and get heads, heads, tails; thus,  $y_1 = 1, y_2 = 1, y_3 = 0$ . There is a true probability of heads called  $\theta$ . We don't know it.

What is the probability of the data? We employ the mass / density function:

$$\begin{aligned}\mathbb{P}(Y_1 = y_1, Y_2 = y_2, Y_3 = y_3; \theta) &= \prod_{i=1}^3 \mathbb{P}(Y_i = y_i; \theta) = \prod_{i=1}^3 \theta^{y_i} (1 - \theta)^{1-y_i} \\ &= \left( \theta^{(1)} (1 - \theta)^{1-(1)} \right) \left( \theta^{(1)} (1 - \theta)^{1-(1)} \right) \left( \theta^{(0)} (1 - \theta)^{1-(0)} \right) \\ &= \theta^2 (1 - \theta)\end{aligned}$$

And now we can calculate the probability of seeing the data assuming  $\theta$ . Assume  $\theta = 0.5$  then,

$$\mathbb{P}(Y_1 = 1, Y_2 = 1, Y_3 = 0; \theta = 0.5) = 0.5^2 (1 - 0.5) = 0.125$$

## The Coin Example from Last Class II

Now we ask the inverse question. If we saw this data  $y_1 = 1, y_2 = 1, y_3 = 0$ , what is the most likely model, i.e. the most likely value of  $\theta$ . We first write down the likelihood function which it's easy because it's the same as the mass / density function

$$\mathcal{L}(\theta; Y_1 = 1, Y_2 = 1, Y_3 = 0) = \mathbb{P}(Y_1 = 1, Y_2 = 1, Y_3 = 0; \theta) = \theta^2(1 - \theta)$$

And now we pick the value of  $\theta$  which maximizes the likelihood,

$$\hat{\theta} := \arg \max_{\theta \in \Theta} \{\mathcal{L}(\theta; \mathbf{x})\}$$

So we need to take the derivative

$$\frac{d}{d\theta} [\theta^2(1 - \theta)] = \frac{d}{d\theta} [\theta^2 - \theta^3] = 2\theta - 3\theta^2$$

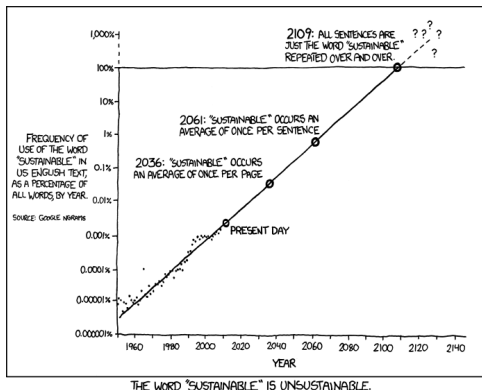
and set it equal to zero:

$$0 = 2\theta - 3\theta^2 = \theta(2 - 3\theta) \Rightarrow 0 = 2 - 3\theta \Rightarrow \hat{\theta}_{\text{MLE}} = \frac{2}{3}$$

i.e. the most likely model for this data is a weighted coin with probability of heads of  $2/3$ .

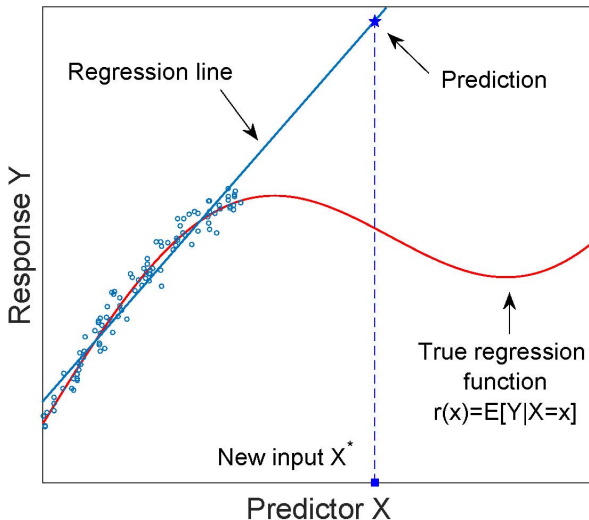
# Extrapolation

Data driven approaches are all focused on accuracy during **interpolation**.



Extrapolation brings trouble. It is important to ask the question for a new observation  $x^*$  if it is within the space of  $x$ 's in the historical data. (Hardly anyone does this... but you should)! Be aware that extrapolation methods of different algorithms differ considerably! [R Demo]

# Reconciliation of those Silly Cartoons



# Dataframe Design

We spoke a lot about featurization i.e. selecting the columns in the dataframe (these are the predictors to measure). Once we did this, we can then go out and sample observations and then measure each for their predictor values.

But we didn't speak at all about selecting the observations themselves (assuming you have some modicum of control of selecting your data). Two things to consider:

- 1 **Generalizability** refers to the ability of the model to generalize, or be **externally valid** when considering new observations. This comes down to sampling observations from the same population as your new data you wish to predict (pretty obvious). Sometimes difficult in practice! Extrapolation??
- 2 Optimal Design

# Optimal Design for Inferring one Slope

Question: assume OLS and that we only care about inference for  $\beta_1$ . We can sample any  $x$  values live in their set  $\mathbb{X}$  e.g.  $\in [x_m, x_M]$ . What should the  $n$  values be?

Let  $x_m = 0$ ,  $x_M = 1$  and  $n = 10$ . The best inference for  $\beta_1$  means ...  $\text{SE}[\hat{\beta}_1]$  is minimum. Design strategies for the  $x$ 's:

- 1 Random sampling
- 2 Uniform spacing:  $\{0, 0.111, 0.222, \dots, 0.999\}$
- 3 Something else?

[R demo]

# Optimal Design: Split Between Extremes

Recall the formula from Stat 102 / 613:

$$\text{SE} [\hat{\beta}_1] = \sqrt{\frac{MSE}{(n-1)s_x^2}}$$

How can we make this small?

- 1 Maximize  $n$  (duh)
- 2 Minimize the numerator,  $MSE$  i.e. minimize the  $SSE$ . Can we do this? Yes by picking the closest  $\hat{\beta}_1$  to  $\beta_1$  (which we already do).
- 3 Maximize the denominator  $(n-1)s_x^2$ . Since  $n$  is already maximized, we can pick  $x_1, \dots, x_n$  to maximize  $s_x^2$ , the sample variance of the predictor. How? Put half of the  $x$ 's at  $x_m$  and the other half at  $x_M$  thereby maximizing the distance from the  $x$ 's to  $\bar{x}$ .



## Optimal Design of Linear Models

We seek the best linear approximation of  $f(x)$  which is  $\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$ . We pick the  $\mathbf{x}$ 's to give us the best linear approximation. What criteria? JMP gives two ways:

❶ Note:  $\text{Var} [\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p] = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$

D-optimality: maximize  $|\mathbf{X}^T \mathbf{X}|$  — this maximizes the variance-covariance among the parameter estimates.

❷ Note:  $\text{Var} [\hat{Y}_1, \dots, \hat{Y}_n] = \sigma^2 \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$

I-optimality: minimize the average prediction variance over the design space.

[R Demo] What did we learn? For linear models with no polynomials or interactions, keep the observations as close to the minimums and maximums as possible. For linear models with polynomials and interactions (more non-parametric than parametric), keep most towards the minimums and maximums and some in the center of the input space.

## Modeling Categorical Responses

Previously the response  $y$  was continuous and via the OLS assumptions we obtained the statistical model,

$$Y \stackrel{\text{ind}}{\sim} \mathcal{N}(\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p, \sigma^2)$$

If the response  $y$  is categorical, can we still use this? No... the only elements in the support of the r.v.  $Y$  are the levels only. [JMP Churn]

First, assume  $Y$  is binary i.e. zero or one. The model (AKA “classifier”) we use is...

$$Y \sim \text{Bernoulli}(f(x_1, \dots, x_p))$$

since  $\mathbb{E}[Y \mid x_1, \dots, x_p] = f(x_1, \dots, x_p)$ , then  $f$  is still the conditional expectation function like before except now it varies only within  $[0, 1]$  and it is the same as  $\mathbb{P}(Y = 1 \mid x_1, \dots, x_p)$ .

## Linear $f(x)$ ?

We can model  $f(x)$  as the simple linear function but this returns values smaller than 0 and larger than 1 and thus it cannot be the conditional expectation function! Why? Lines vary between  $(-\infty, +\infty)$ .

We need a “link function” to connect the linear function to the restricted support of the response:

$$\lambda(f_{\mathbb{R}}(x_1, \dots, x_p)) = f(x_1, \dots, x_p)$$

And the parametric assumption would be

$$\lambda(s_{\mathbb{R}}(x_1, \dots, x_p; \theta_1, \dots, \theta_\ell)) = s(x_1, \dots, x_p; \theta_1, \dots, \theta_\ell)$$

And assuming a linear form of  $s_{\mathbb{R}}$ ,

$$\lambda(\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p) = ?$$

## Choice of $\lambda$ ?

We just need  $\lambda : \mathbb{R} \rightarrow [0, 1]$ . There are infinite  $\lambda$ 's to choose from. I've only seen three used:

- 1 Logistic link:  $\lambda(w) = \frac{e^w}{1+e^w}$  (most common)
- 2 Inverse normal (probit) link:  $\lambda(w) = \Phi^{-1}(w)$  where  $\Phi$  is the normal CDF function (somewhat common)
- 3 Complementary Log-log (cloglog) link:  $\lambda(w) = \ln(-\ln(w))$  (rare!)

Let's investigate what the first one means. Define  $p := \mathbb{P}(Y = 1)$ . We can think about probability in another way:

$$\text{odds}(Y = 1) := \frac{p}{1 - p}$$

So if odds = 4:1, what is  $p$ ? This means that the probability of the event happening is four times more likely than the complement happening. Or... of 4+1 runs, 4 will be a yes. What is the range of odds?  $[0, \infty)$ .

# Why Logistic Link is Interpretable

Now let's take the log odds (called the logit function):

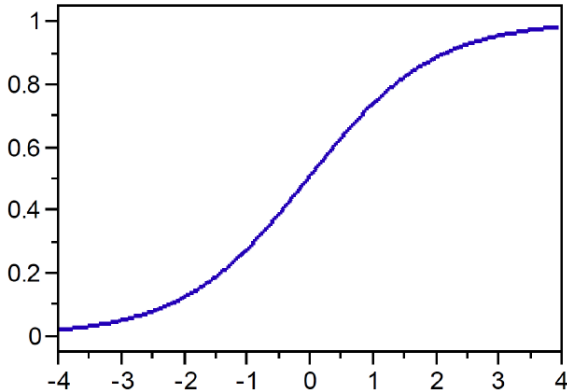
$$\text{logit}(Y = 1) := \ln(\text{odds}(Y = 1)) = \ln\left(\frac{p}{1-p}\right)$$

What is the range of the logit function? All of  $\mathbb{R}$ . Hence, we can now set this equal to our  $s_{\mathbb{R}}$  function. In the linear modeling context,

$$\begin{aligned}\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p &= \text{logit}(Y = 1) = \ln\left(\frac{p}{1-p}\right) \\ e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p} &= \frac{p}{1-p} \\ (1-p)e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p} &= p \\ e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p} &= p + pe^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p} \\ p &= \frac{e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p}}{1 + e^{\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p}} = \lambda(\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p)\end{aligned}$$

Thus, a change in the linear model becomes a linear change in log-odds. This is (I would say) the most interpretable link function situation we've got.

# The Logistic Function



## How to Obtain a Model Fit

A model fit would mean we estimate  $\{\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p\}$ . We initially did this estimation for regression (continuous  $y$ ) by defining a loss function, SSE, and finding the optimal solution via calculus. What do we do now??

Likelihood to the rescue. First the “logistic regression assumptions”

- 1 Linear-Logistic conditional expectation
- 2 Independence

$$\begin{aligned} & \mathbb{P}(Y_1 = y_1, Y_2 = y_2, \dots, Y_n = y_n \mid \mathbf{X}_1 = \mathbf{x}_1, \mathbf{X}_2 = \mathbf{x}_2, \dots, \mathbf{X}_n = \mathbf{x}_n) \\ &= \prod_{i=1}^n \mathbb{P}(Y_i = y_i \mid \mathbf{X}_1 = \mathbf{x}_i) \end{aligned}$$

How?

# Maximum Likelihood Estimates

$$= \prod_{i=1}^n p^{y_i} (1 - p)^{1-y_i}$$

How?

$$\begin{aligned} & \mathcal{L}(\beta_0, \beta_1, \dots, \beta_p; \mathbf{x}_1, \dots, \mathbf{x}_n) \\ &= \prod_{i=1}^n \left( \frac{e^{\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}}}{1 + e^{\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}}} \right)^{y_i} \left( 1 - \frac{e^{\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}}}{1 + e^{\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}}} \right)^{1-y_i} \end{aligned}$$

How? This does not have a simple, closed form solution. The computer iterates numerically using gradient methods. It usually uses the  $\ln(\cdot)$  of above, since it's (1) numerically more stable and (2) the expression is easier to work with. When it “converges” on the values of the parameters that maximize the above, these are shipped to you as  $\{\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p\}$ . This is called “running a logistic regression”. This usually is instant on a modern computer.



# Prediction with Logistic Regression

How?

$$\hat{p} = \hat{p}(x_1^*, \dots, x_p^*) = \frac{e^{\hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p}}{1 + e^{\hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p}}$$

Note the predictions are for the conditional expectation function, the probability itself, the **estimated expected probability**. However, you may actually wish to predict the response, the 1 or the 0. What to do?

You can create a **classification rule** which allows you to make a decision about the response based on the probability. What is the most intuitive classification rule?

$$\hat{y} = \mathbb{1}_{\hat{p} \geq 0.5} := \begin{cases} 1 & \text{if } \hat{p} \geq 0.5 \\ 0 & \text{if } \hat{p} < 0.5 \end{cases}$$

AKA the “most likely criterion”. We will return to prediction and evaluation of predictive performance later but first... inference.

# Global Test in Logistic Regression

Recall in OLS regression, gaussian (normal) theory directly gave us  $t$ -tests and  $F$ -tests. Under the logistic regression assumptions, **we have no such analogous theory!** However, we can make use of the ... likelihood ratio test. Recall:

$$LR := \max_{\theta \in \Theta} \mathcal{L}(\theta; \mathbf{x}) / \max_{\theta \in \Theta_R} \mathcal{L}(\theta; \mathbf{x})$$

Let's now do a "whole model" / "global" / "omnibus" test:

$$H_0 : \beta_1 = 0, \beta_2 = 0, \dots, \beta_p = 0, \quad H_a : \text{at least one is non-zero}$$

So  $\Theta$  would be the space of all  $\beta_0, \beta_1, \dots, \beta_p$  and  $\Theta_R$  will restrict the space to only  $\beta_0$  with zeroes for all other "slope" parameters.

$$LR = \frac{\max_{\beta_0, \beta_1, \dots, \beta_p} \mathcal{L}(\beta_0, \beta_1, \dots, \beta_p; y_1, \dots, y_n, \mathbf{x}_1, \dots, \mathbf{x}_n)}{\max_{\beta_0} \mathcal{L}(\beta_0, \beta_1 = 0, \dots, \beta_p = 0; y_1, \dots, y_n, \mathbf{x}_1, \dots, \mathbf{x}_n)}$$

So on top the computer iterates to find  $\{\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p\}$ , plugs it in and computes the likelihood and on the bottom the computer independently iterates to find  $\{\hat{\beta}_0\}$ , plugs it in and computes the likelihood, then together, the  $LR$ .

## Partial Tests in Logistic Regression

We then look at  $Q = 2 \ln(LR)$  and compare it to the appropriate  $\chi^2$  distribution. Here, since we've dropped  $p$  parameters / degrees of freedom, we look at the critical  $\chi^2_{p,\alpha}$  value.

Let's say we want to test something like:

$$H_0 : \beta_1 = 0 \ \& \ \beta_2 = 0, \quad H_a : \text{at least one is non-zero}$$

We can again use the likelihood ratio test:

$$LR = \frac{\max_{\beta_0, \beta_1, \dots, \beta_p} \mathcal{L}(\beta_0, \beta_1, \dots, \beta_p; y_1, \dots, y_n, \mathbf{x}_1, \dots, \mathbf{x}_n)}{\max_{\beta_0, \beta_3, \dots, \beta_p} \mathcal{L}(\beta_0, \beta_1 = 0, \beta_2 = 0, \beta_3, \dots, \beta_p = 0; y_1, \dots, y_n, \mathbf{x}_1, \dots, \mathbf{x}_n)}$$

We then look at  $Q = 2 \ln(LR)$  and compare it to the appropriate  $\chi^2$  distribution. Here, since we've dropped 2 parameters / degrees of freedom, we look at the critical  $\chi^2_{2,\alpha}$  value.

# Individual Tests in Logistic Regression

Let's say we want to test an individual slope coefficient:

$$H_0 : \beta_j = 0, \quad H_a : \beta_j \neq 0$$

We can again use the likelihood ratio test:

$$LR = \frac{\max_{\beta_0, \beta_1, \dots, \beta_p} \mathcal{L}(\beta_0, \beta_1, \dots, \beta_p; y_1, \dots, y_n, \mathbf{x}_1, \dots, \mathbf{x}_n)}{\max_{\beta_0, \beta_1, \dots, \beta_{j-1}, \beta_{j+1}, \dots, \beta_p} \mathcal{L}(\beta_0, \beta_1, \dots, \beta_{j-1}, \beta_j = 0, \beta_{j+1}, \dots, \beta_p; y_1, \dots, y_n, \mathbf{x}_1, \dots, \mathbf{x}_n)}$$

We then look at  $Q = 2 \ln(LR)$  and compare it to the appropriate  $\chi^2$  distribution. Here, since we've dropped 1 parameter / degrees of freedom; thus we look at the critical  $\chi^2_{1, \alpha}$  value.

There is something special about a  $\chi^2$  r.v. with one degree of freedom. Note this cool fact from probability theory:  $\sqrt{\chi^2_1} \sim \mathcal{N}(0, 1)$  i.e. a “z-score”. This is how JMP produces standard errors for logistic regression coefficients.

# Telecom Company Churn Example

In marketing “churn” refers to a customer canceling their service. Studies suggest that it costs 5-10x more to acquire a new customer than to retain an old customer. Thus, predicting churn is of major interest!

Here's a dataset from a telecom company (likely it's churn on Verizon / AT&T / T-Mobile / Sprint's cell-phone plan). We have 7,043 customers with 20 features. This is likely a nearly-mindless dump!! Churn is defined to be a complete cancellation of services in the next month period. Since we are predicting churn, define  $y = 1$  to be churn, so the  $\hat{p}$ 's are estimates of probability of churning (this just makes everything easier to understand).

We begin just trying to model  $y$ : churn vs.  $x$ : tenure (the number of months customer is currently subscribed for). What do you think the relationship will be?  $\partial/\partial x[f(x)]$ ?

## Results of Simple Logistic Regression

[JMP demo] Which likelihood numbers are best? Well highest likelihoods are good which means highest log likelihoods which means *smallest* negative log likelihoods.... yup it's confusing. Here's what JMP did:

$$\begin{aligned} Q &= 2 \ln(LR) = 2 \left( \ell(\hat{\theta}; x) - \ell(\hat{\theta}_R; x) \right) \\ &= 2 \left( -(3595.9341) - -(4075.0729) \right) \\ &= 2(479.1389) = 958.2778 \\ &= 2 \ln(1.22 \times 10^{208}) \end{aligned}$$

and  $\chi^2_{1,5\%} = \text{qchisq}(.95, 1) = 3.84$ . So this passes the test (comfortably). We reject  $H_0$  and conclude that the model is linearly useful. ALSO: equivalent to a test of one variable. We reject  $H_0$  and conclude tenure has a linear effect on the log-odds of churn.

## Basic Predictions I

Predict estimated expected probability of churn for someone who has 1 month of tenure

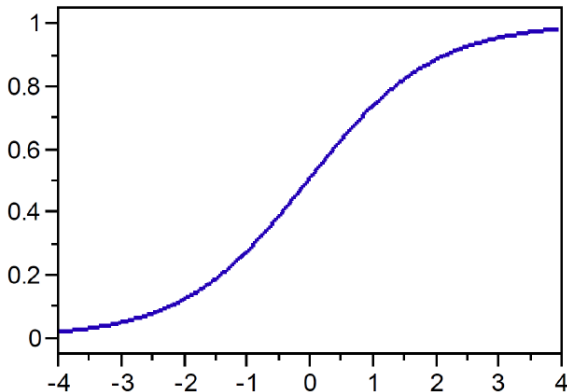
$$\hat{p} = \frac{e^{\hat{\beta}_0 + \hat{\beta}_1 x_1}}{1 + e^{\hat{\beta}_0 + \hat{\beta}_1 x_1}} = \frac{e^{0.0273 + (-0.0288)(1)}}{1 + e^{0.0273 + (-0.0288)(1)}} = \frac{0.9985}{1.9985} = 0.500$$

How about 2 months of tenure?

$$\hat{p} = \frac{e^{\hat{\beta}_0 + \hat{\beta}_1 x_1}}{1 + e^{\hat{\beta}_0 + \hat{\beta}_1 x_1}} = \frac{e^{0.0273 + (-0.0288)(2)}}{1 + e^{0.0273 + (-0.0288)(2)}} = \frac{0.9702}{1.9702} = 0.492$$

i.e. a difference in about 0.8% as measured on a probability scale.

# The Logistic Function



A move of one unit in  $x$  when  $x \approx 0$  is a much bigger move than one unit in  $x$  when  $x \approx 3$



## Parameter Standard Error

To add to the confusion... JMP prefers to calculate parameter estimates and standard error via the Wald test, which is similar to the likelihood ratio test. Thus,  $761.00 \neq 958.28$  but, remarkably, they are about the same conceptually – both large and significant. The standard errors,

$$\underbrace{Q = z^2}_{\substack{\text{fact from} \\ \text{probability} \\ \text{theory}}} = \left( \frac{\hat{\beta}_1}{s_{\hat{\beta}_1}} \right)^2 \Rightarrow s_{\hat{\beta}_1} = \frac{|\hat{\beta}_1|}{\sqrt{Q}}$$

are expectedly about the same

$$s_{\hat{\beta}_1} = \frac{|-0.0387682|}{\sqrt{761.00}} = 0.0014 \quad (\text{via the Wald test})$$

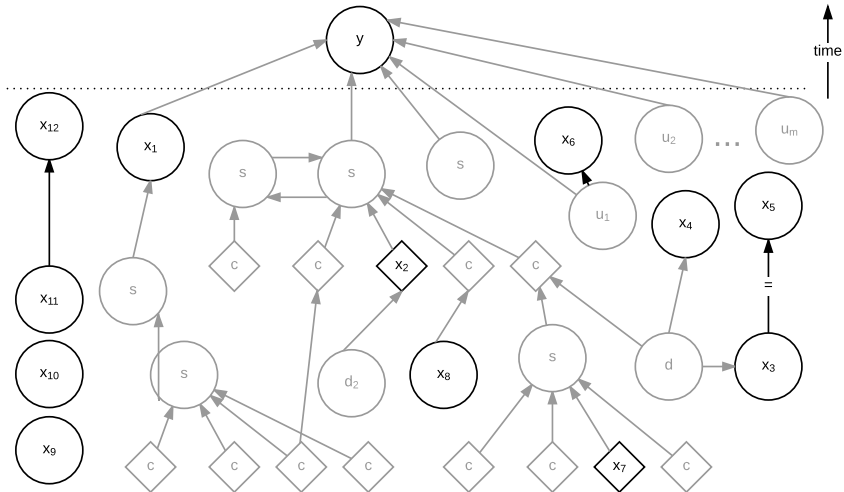
$$s_{\hat{\beta}_1} = \frac{|-0.0387682|}{\sqrt{958.28}} = 0.0013 \quad (\text{via the LR test})$$

# Multivariate Logistic Regr. Interp. I

Now let's use all variables. Questions:

- Which variable(s) should we leave out? `customerID` — no causal mechanism.
- Why are we getting biased estimates and zeroes? Perfect multicollinearity. Solution? Kill columns until you don't get the error anymore. Turns out 6 have to be removed.

# Realistic Predictors Illustration (updated)

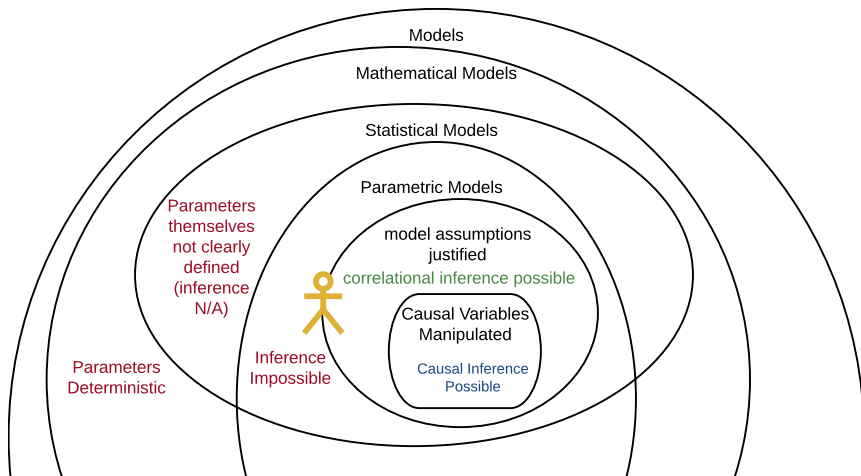


# Multivariate Logistic Regr. Interp. II

Now let's use all variables. Questions:

- Which variable(s) should we leave out? `customerID` — no causal mechanism.
- Why are we getting biased estimates and zeroes? Perfect multicollinearity.
- Are all these variables “significant”? No. After a Bonferroni correction, we “lose” two. Are they for sure insignificant? No, maybe we didn't have enough power to detect them...
- Do the coefficients sign / size make sense intuitively? Yes?
- Which of the “driving” variables are *not* able to be manipulated? Senior citizen, total charges (?)
- Should you (the manager) try to incentivize electronic checks? Paperless billing? Dropping multiple phone lines? You can try these things if you have nothing to lose, but remember, they are not guaranteed to be causal! And they may **backfire!!** (Can you think of an example??)

# Remember Where you At!!



# Evaluating Logistic Regression Models

Many, many measures reported by JMP that we generally don't use:

- $R^2(U)$
- Generalized  $R^2$
- Mean Negative Log probability
- "RMSE"
- Mean Absolute Deviation

And ones that we do use:

- AICc / BIC (for something a little bit different... we will come back to this in a couple of lectures)
- Misclassification Rate

We now cover evaluating classification models in general (not only in the context of logistic regression models specifically).

## Probability Predictions $\Rightarrow$ Level Predictions

Recall that ... you can create a classification rule which allows you to make a decision about the response based on the probability. The most intuitive classification rule is:

$$\hat{y} = \mathbb{1}_{\hat{p} \geq 0.5} := \begin{cases} 1 & \text{if } \hat{p} \geq 0.5 \\ 0 & \text{if } \hat{p} < 0.5 \end{cases}$$

In regression, you examined functions of the residuals  $e_i := y_i - \hat{y}_i$  to assess model fit. What is an analogous residual here? There are four residuals, two representing errors. The best way to see them is to create the **confusion matrix**:

		$\hat{y}$ (decision)	
		1	0
$y$ (truth)	1	true positive (TP)	false negative (FN)
	0	false positive (FP)	true negative (TN)

Why do “correlations rock” here?? We are purely evaluating predictive performance... no inferential claims!

## Confusion Matrix for Churn Model

JMP gives us the matrix [JMP], but they don't annotate it well.  
Here are some numbers I like to see:

		$\hat{y}$		Totals	Model Errors
		1	0		
$y$	1	$TP = 1012$	$FN = 857$	$P = 1869$	$FNR = 45.9\%$
	0	$FP = 531$	$TN = 4632$	$N = 5163$	$FPR = 10.2\%$
Totals		$\hat{P} = 1543$	$\hat{N} = 5489$	$n = 7032$	
Use errors		$FDR = 34.3\%$	$FOR = 15.6\%$		$ME = 19.7\%$

There are other metrics commonly reported

- Sensitivity =  $\frac{TP}{TP+FN} = \frac{TP}{P}$ , the proportion of positives successfully recovered (large value = good model), 54.9% above
- Specificity =  $\frac{TN}{TN+FP} = \frac{TN}{N}$ , the proportion of negatives successfully recovered (large value = good model), 89.8% above



# Misclassification Error

Already... what is one broad, general metric to evaluate the model?  
Misclassification error cost function (or Accuracy):

$$ME := \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{y_i \neq \hat{y}_i}$$

$$ACC := 1 - ME = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{y_i = \hat{y}_i}$$

This essentially treats both types of errors (FN and FP) equally (more on this later).

# There's are a Ton of Metrics...

From wikipedia...

		Predicted condition			
Total population		Predicted Condition positive	Predicted Condition negative	Prevalence = $\frac{\sum \text{Condition positive}}{\sum \text{Total population}}$	
True condition	condition positive	True positive	False Negative (Type II error)	True positive rate (TPR), Sensitivity, Recall, probability of detection = $\frac{\sum \text{True positive}}{\sum \text{Condition positive}}$	False negative rate (FNR), Miss rate = $\frac{\sum \text{False negative}}{\sum \text{Condition positive}}$
	condition negative	False Positive (Type I error)	True negative	False positive rate (FPR), Fall-out, probability of false alarm = $\frac{\sum \text{False positive}}{\sum \text{Condition negative}}$	True negative rate (TNR), Specificity (SPC) = $\frac{\sum \text{True negative}}{\sum \text{Condition negative}}$
Accuracy (ACC) = $\frac{\sum \text{True positive} + \sum \text{True negative}}{\sum \text{Total population}}$		Positive predictive value (PPV), Precision = $\frac{\sum \text{True positive}}{\sum \text{Test outcome positive}}$	False omission rate (FOR) = $\frac{\sum \text{False negative}}{\sum \text{Test outcome negative}}$	Positive likelihood ratio (LR+) = $\frac{\text{TPR}}{\text{FPR}}$	Diagnostic odds ratio (DOR) = $\frac{\text{LR+}}{\text{LR-}}$
		False discovery rate (FDR) = $\frac{\sum \text{False positive}}{\sum \text{Test outcome positive}}$	Negative predictive value (NPV) = $\frac{\sum \text{True negative}}{\sum \text{Test outcome negative}}$	Negative likelihood ratio (LR-) = $\frac{\text{FNR}}{\text{TNR}}$	

Others I will be talking about:

- False Discovery Rate (FDR) =  $\frac{FP}{TP+FP} = \frac{FP}{\hat{P}}$ , the proportion of positives of those predicted to be positive (small value = good model)
- False Omission Rate (FOR) =  $\frac{FN}{TN+FN} = \frac{FN}{\hat{N}}$ , the proportion of negatives of those predicted to be negative (small value = good model)

## Generalizing the Classification Rule

Recall the classification rule  $\hat{y} = \mathbb{1}_{\hat{p} \geq 0.5}$ . Using 0.5 is a principled default but we can use any rule  $p_0 \in (0, 1)$ :

$$\hat{y} = \mathbb{1}_{\hat{p} \geq p_0} := \begin{cases} 1 & \text{if } \hat{p} \geq p_0 \\ 0 & \text{if } \hat{p} < p_0 \end{cases}$$

What happens when we change the  $p_0$  threshold? If  $p_0 \uparrow \Rightarrow \hat{P} \downarrow$  and  $\hat{N} \uparrow$ . If  $p_0 \downarrow \Rightarrow \hat{P} \uparrow$  and  $\hat{N} \downarrow$ . Changing  $p_0$  changes the column totals and obviously creates a whole new confusion matrix.

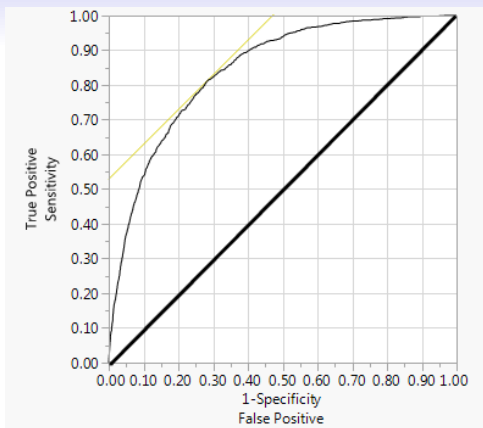
So now it's simple, vary  $p_0$  and pick the best model according to your cost / error / loss function (the *ME* at the moment). Let's just do every  $p_0$ !

# Receiver-Operator Characteristic Table

ROC Table							
Prob	1-Specificity	Sensitivity	Sens- (1-Spec)	True Pos	True Neg	False Pos	False Neg
.	0.0000	0.0000	0.0000	0	5163	0	1869
0.8117	0.0000	0.0005	0.0005	1	5163	0	1868
0.8104	0.0000	0.0011	0.0011	2	5163	0	1867
0.8093	0.0000	0.0016	0.0016	3	5163	0	1866
0.8092	0.0000	0.0021	0.0021	4	5163	0	1865
0.8090	0.0000	0.0027	0.0027	5	5163	0	1864
0.8085	0.0000	0.0032	0.0032	6	5163	0	1863
0.8083	0.0000	0.0037	0.0037	7	5163	0	1862
0.8082	0.0000	0.0043	0.0043	8	5163	0	1861
0.8079	0.0000	0.0048	0.0048	9	5163	0	1860
0.8079	0.0000	0.0054	0.0054	10	5163	0	1859
0.8077	0.0000	0.0059	0.0059	11	5163	0	1858
0.8076	0.0002	0.0059	0.0057	11	5162	1	1858
0.8072	0.0002	0.0064	0.0062	12	5162	1	1857
0.8065	0.0002	0.0070	0.0068	13	5162	1	1856
0.8064	0.0002	0.0075	0.0073	14	5162	1	1855
0.8061	0.0002	0.0080	0.0078	15	5162	1	1854

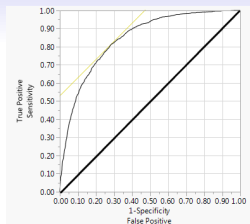
Here, Prob is what we denoted  $p_0$ . “Best model” is not defined here by highest *ACC* (lowest *ME*), it’s determined by highest **specificity + sensitivity** or equivalently, the highest **sensitivity - (1 - specificity)**. JMP indicates that row with a ★. This is an arbitrary metric, but is a good default.

# Receiver-Operator Characteristic Curve



This is graphical illustration of the table. Each dot represents the sensitivity-specificity tradeoff for each  $p_0$ . The starred row of maximum sensitivity + specificity is indicated here by a yellow tangent line. I drew the diagonal line to indicate predictive performance that is expected “by chance”. Why?

# Area Under the Curve (AUC) Metric



If you built a model by chance the “area under the curve” (or to the right of the curve) on the graph would be ... 0.5 since the graph is a unit square. Under the ROC curve itself (or to its right) is an area ... greater than 0.5. Here, it's 0.844. This metric is called AUC and is widely used as a metric to assess performance of all possible classifiers in this set of models together, it is a composite metric unlike *ME* or anything derived from an individual confusion table.

AUC is nice to evaluate overall performance of all possible models... but at the end of the day... you ship **ONE** model! So we still need a means of evaluating our one model from one confusion table.

## Churn Example Where $p_0 = 0.10$

$p_0 = 0.5$		$\hat{y}$		Totals	Model Errors
		1	0		
$y$	1	$TP = 1012$	$FN = 857$	$P = 1869$	$FNR = 45.9\%$
	0	$FP = 531$	$TN = 4632$	$N = 5163$	$FPR = 10.2\%$
Totals		$\hat{P} = 1543$	$\hat{N} = 5489$	$n = 7032$	
Use errors		$FDR = 34.3\%$	$FOR = 15.6\%$		$ME = 19.7\%$

$p_0 = 0.1$		$\hat{y}$		Totals	Model Errors
		1	0		
$y$	1	$TP = 1772$	$FN = 97$	$P = 1869$	$FNR = 5.1\%$
	0	$FP = 2669$	$TN = 2494$	$N = 5163$	$FPR = 51.6\%$
Totals		$\hat{P} = 4441$	$\hat{N} = 2591$	$n = 7032$	
Use errors		$FDR = 60.1\%$	$FOR = 3.7\%$		$ME = 39.3\%$

Which numbers did not change?  $n$ ,  $P$  and  $N$ . Why? These are fixed according to the dataframe. All other numbers changed! What happened to our first means of evaluation, the Misclassification Error? It increased from  $19.7\% \rightarrow 39.3\%$ . So isn't this a worse model??

Not necessarily... It depends on what your goal is!

## Asymmetric Costs in a Classifier

These are always two types of errors but the costs are not always the same.

$p_0 = 0.1$		$\hat{y}$		Totals	Model Errors
		1	0		
$y$	1	$TP = 1772$	$FN = 97$	$P = 1869$	$FNR = 5.1\%$
	0	$FP = 2669$	$TN = 2494$	$N = 5163$	$FPR = 51.6\%$
Totals		$\tilde{P} = 4441$	$\tilde{N} = 2591$	$n = 7032$	
Use errors		$FDR = 60.1\%$	$FOR = 3.7\%$		$ME = 39.3\%$

Imagine we really are the Telecom business manager. It costs 5-10x more to acquire a new customer than to engage a customer who is likely to churn. What type of error specifically is *very* costly? The *FN*. Who are they? These are those who you said were not going to churn *and they did!* Cost? You need to acquire a new customer! The other type of error is less costly, the *FP*. Who are they? These are the people you thought were going to churn and did not. Cost? Whatever the incentive package is.



## Weighted Misclassification Error

We now define two costs: (1) the cost of the *FP* denoted  $c_{FP}$  and (2) the cost of the *FN* denoted  $c_{FN}$ . We then define the weighted misclassification error evaluation metric:

$$ME_w := \frac{1}{n} \sum_{i=1}^n c_{FP} \mathbb{1}_{y_i=0 \& \hat{y}_i=1} + c_{FN} \mathbb{1}_{y_i=1 \& \hat{y}_i=0}$$

We now vary  $p_0$  to locate the model that optimizes this error to be minimum.

# Minimum Weighted Misclassification Error

Let's assume that  $c_{FN} = \$1000$  and  $c_{FP} = \$100$  just for the example's sake. Note: this is a **cost ratio** of 10:1.

	Prob	TP	TN	FP	FN	COST
1	0.8117	1	5163	0	1868	1868000
2	0.8104	2	5163	0	1867	1867000
3	0.8093	3	5163	0	1866	1866000
4	0.8092	4	5163	0	1865	1865000
5	0.8090	5	5163	0	1864	1864000
6	0.8085	6	5163	0	1863	1863000
7	0.8083	7	5163	0	1862	1862000
8	0.8082	8	5163	0	1861	1861000
9	0.8079	9	5163	0	1860	1860000

We now calculate the cost and find the minimum model (i.e. the  $p_0$  to ship). [JMP] Or alternatively, we can select the model with the closest  $FN/FP \approx 10 : 1$  to match the stakeholder preference of the desired cost ratio.

## $\hat{p}$ 's as Ordinal Values

One final point... If we were on a mission to find the top  $m$  churners. What would we do? Sort the  $\hat{p}$ 's and return the top  $m$ .