# Regression Models

## Role of Regression Models in Clinical Research:

The practical goal of epidemiology is to measure and interpret associations between suspected risk factors and outcomes. For causal research, the usual measurement goal of the investigation is to quantifying the effect of a single risk factor of interest (exposure) on the outcome while controlling for confounding by other factors (**explanation**). However, a second measurement goal of epidemiology (especially of clinical epidemiology) may be to quantify the joint effect of many risk factors to estimate an individual's risk of developing or possessing an outcome (**prediction**). Regression Models are commonly used to achieve either of these goals. However, the steps used to develop these models and the focus on their results depends on whether explanation or prediction is the goal of the analysis.

## Regression Coefficients

Regression coefficients can be interpreted as slope coefficients, reflecting the change in an outcome, per unit chance in a risk factor. For example, the following data show the relationship between smoking (measured in packs smoked per day) and the log(odds of dying), the **logit**, during the 24 years of follow-up in the Framingham Heart Study (FHS) Teaching Data Set. (N.B. 32 of the 4434 participants have missing values for cigpday1 and therefore have missing values for packs smoked per day.)

| Packs Smoked Per Day | Number at Risk | Number of Deaths | Estimated Risk | Logit |
|---|---|---|---|---|
| 0 | 2253 | 762 | 762/2253 = .34 | log(.34/.66) = -.66 |
| 1 | 1671 | 573 | 573/1671 = .34 | log(.34/.66) = -.66 |
| 2 | 398 | 169 | 169/398 = .42 | log(.42/.58) = -.32 |
| 3+ | 80 | 36 | 36/80 = .45 | log(.45/.55) = -.20 |

The data in this table shows evidence of an increasing log(odds of death) with increasing number of packs of cigarettes smoked, which might be approximated by the following linear equation

$$\log(\text{Odds of Death}) = B_0 + B_1(\text{Packs})$$

A regression equation that describes the relationship between the log(odds of an outcome) as function of one or more risk factors is called a **Logistic Regression Model**. The slope of this linear equation ($B_1$) measures the change in the log(Odds of Death) per smoking one additional pack of cigarettes per day. For example, if $P_x$ is the risk of death when smoking "x" packs of cigarettes per day then

$$B_1 = \log(P_{x+1}/(1-P_{x+1})) - \log(P_x/(1-P_x))$$

$$= \log[(P_{x+1}/(1-P_{x+1})) / \log(P_x/(1-P_x))]$$

$$= \log \text{ (Odds Ratio)}$$

The last equation shows that regression coefficients for a logistic regression model can be interpreted as the logarithm of a common measure of association, the Odds Ratio. In addition if the logistic regression model contains multiple risk factors, the coefficient for any risk factor has the interpretation as the log(Odds Ratio), measuring the association between that risk factor and the outcome, conditional on all other risk factors in the model.

For example, the following formula describes the risk (P) of death during 24 years of follow-up in the FHS Teaching Data Set as a function of five risk factors: current smoking status (CURSMOK1: 1=yes, 0=no), age in years (AGE), male sex (MALE: 1= yes, 0=no), hypertension (HIGHBP1: 1 if sysbp1 $\geq$ 140 or diabp1 $\geq$ 90, 0 otherwise) and diabetes (DIABETES1: 1=yes, 0=no).

$$\log(P/(1-P)) = B_0 + B_1\text{CURSMOK1} + B_2\text{AGE} + B_3\text{MALE} + B_4\text{HIGHBP1} + B_5\text{DIABETES1}$$

When fit to the FHS teaching data set, the fitted model becomes

$$\log(P/(1-P)) = -7.5869 + 0.5522(\text{CURSMOK1}) + 0.1181(\text{AGE}) + 0.7759(\text{MALE})$$
$$+ 0.6386(\text{HIGHBP1}) + 1.5834(\text{DIABETES1})$$

If this model correctly describes the relationship between the 5 risk factors and the risk of death, then it follows that the 24-year risk of death for a 50-year-old female, without hypertension and without diabetes is

$$\log(P/(1-P)) = -7.5869 + 0.5522(\text{CURSMOK1}) + 0.1181(50) + 0.7759(0) + 0.6386(0)$$
$$+ 1.5834(0)$$

$$= -1.6819 + 0.5522(\text{CURSMOK1})$$

This equation resembles the linear equation presented above. Therefore, the coefficient of CURSMOK1 (0.5522) is the log(Odds Ratio) describing the relationship between smoking and death, for a 50-year-old female, without hypertension and without diabetes and

$$OR = e^{0.5522} = 1.74$$

On the other hand, the 24-year risk of death for a male, 50-year-old male with hypertension and with diabetes is

$\log(P/(1-P))$    = -7.5869 + 0.5522(CURSMOK1) + 0.1181(50) + 0.7759(1) + 0.6386(1)
            + 1.5834(1)

        = 1.3160 + 0.5522(CURSMOK1)

The coefficient of CURSMOK1 (0.5522) again is the log(Odds Ratio) describing the relationship between smoking and death, but now for a 50-year-old male, with hypertension and without diabetes and

$$OR = e^{0.5522} = 1.74$$

In summary, for any combination of (AGE, MALE, HIGHBP1, and DIABETES1), the model's estimate of the effect of CURSMOKE1 is OR = 1.74.

In general, a regression coefficient estimates the effect of a predictor on the outcome, controlling for all other factors in the model. This property is the basis for controlling for confounding by a regression model requires. This method requires including sufficient terms in the model to represent the confounders. It also requires valid **modeling assumptions**, and wrong modeling assumptions can lead to wrong conclusions.

**Multiple Regression Models**

A multiple regression model is a mathematical expression that postulates a relationship between an outcome and a set of predictors. Typically the predictors in the model represent the exposure of interest, potential confounders, and possibly effect modifiers. Perhaps the most commonly used model in epidemiologic research is the **logistic regression model**. This model is appropriate for a **binary outcome** (e.g. dead versus alive) and describes the risk (P) of developing the outcome as a function of the predictors $(X_1, X_2, \ldots, X_n)$ by the following formula:

$$\log(P/(1-P)) = B_0 + B_1X_1 + B_2X_2 + \ldots + B_nX_n$$

The unknown coefficients in this model are the intercept term ($B_0$) and the coefficients

($B_i$) of the predictors ($X_i$). The intercept term ($B_0$) specifies the value for the outcome ($\log(P/(1-P))$) when all predictors are set equal to zero (i.e. $X_i = 0$, $i = 1,2, \ldots, n$). More importantly, each coefficient ($B_i$) is a slope coefficient, describing the change in $\log(P/(1-P))$ per unit change in $X_i$ ($\log(OR)$) when all other predictors are fixed. This coefficient is often interpreted as an estimate of the "effect" of the corresponding predictor, controlling for all of the other factors in the model. In this way multivariate models are the most common method for measuring the effect of an intervention, controlling for confounding by other factors.

Another multivariable model often used in clinical research is the **linear regression model**, which describes the relationship between a **continuous outcome** (Y) and a set of predictors. This model assumes that the average outcome (expected value, E(Y)) is related to the predictors according to the following formula

$$E(Y) = B_0 + B_1X_1 + B_2X_2 + \ldots + B_nX_n$$

Other models that are sometimes used in clinical research are the **Cox Regression Model** for survival time outcomes with censoring and the **Poisson Regression Model** for count outcomes. The main differences in these models are the type of outcome (binary outcome, continuous outcome, survival time outcome, and count outcome) and the method for fitting a model to a data set. However, all model share many principles in common, including the interpretation of any regression coefficient as representing the change in the outcome per unit change in a predictor, controlling for all other predictors in the model.

Methods for estimating regression coefficients will not be discussed in detail in these lecture notes. Algorithms for fitting a model to a data set differ, depending on the type of the regression model. For example, maximum likelihood estimation is used to estimate the coefficients of a logistic regression model. Least estimation is used to estimate the coefficients in a linear regression model. Partial likelihood methods are used to estimate the coefficients in a Cox regression model. All methods have the goal that outcome estimates for individuals from the model should agree as much as possible with the actual outcomes. For example, risk estimates from a logistic regression model should be as high as possible (close to 1.0) for those who develop the outcome and as low as possible (close to 0.0) for those who do not develop the outcome.

**Model Assumptions**

The following table shows the equation for the Mortality Prediction Model ($MPM_0$) (*Lemeshow et al. Predicting the Outcome of Intensive Care Unit Patients. J Am Stat Assoc 1988;83(402):348-356.*). The $MPM_0$ was developed from a logistic regression model to predict the risk of in-hospital death, P, from a set of clinical measures for patients admitted to an ICU.

**Table**. Model predicting the risk of hospital mortality (P) for patients admitted to an intensive care unit based on a logistic regression model containing seven predictors

Predictors

        CONS        : level of consciousness   (1 if coma or deep stupor, 0 otherwise)
        TYPE        : type of admission (1 if emergent, 0 if elective)
        CANCER   : cancer as part of present problem (1 if yes, 0 if no)
        CPR         : prior CPR (1 if yes, 0 if no)
        INFECT    : infection  (1 if probable, 0 otherwise)
        AGE         : age in ten year increments
        SBP         : systolic blood pressure
        SBP2       : SBP squared.

Model :

$$\log(P/(1-P)) = -1.370 + 2.44(CONS) + 1.81(TYPE) + 1.49(CANCER) + .974(CPR)$$
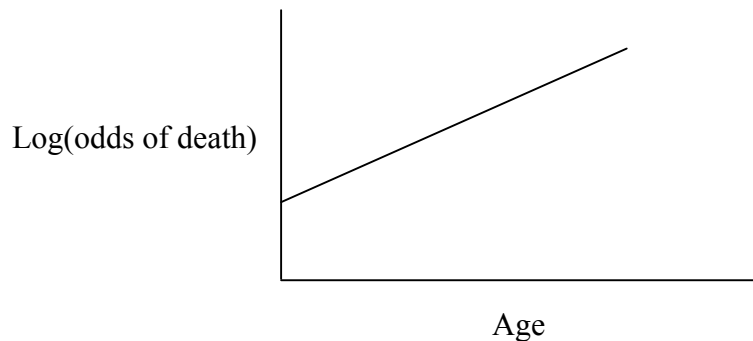$$+ .965(INFECT) + .0368(AGE) - .0606(SBP) + .000175(SBP2)$$

Most of the factors in the MPM model are binary predictors, each representing the presence or absence of a risk factor. However, the model also contains terms for two continuous risk factors: age and systolic blood pressure. The model assumes that each predictor has a single effect on the outcome as measured by its coefficient and that this effect holds over all subgroups of subjects that are defined by the other predictors in the model. For example, this model assumes that the effect prior CPR (as estimated by its regression coefficient (0.974) is not modified by any of the other predictors in the model. This condition is known as the **assumption of additivity**.

The model in this table also contains a single term for the age of each subject. If we fix the values for the other predictors, then this model assumes that the conditional linear relationship between age and the log odds of dying, which is described by the following equation and represented by a straight line as depicted in the following figure.
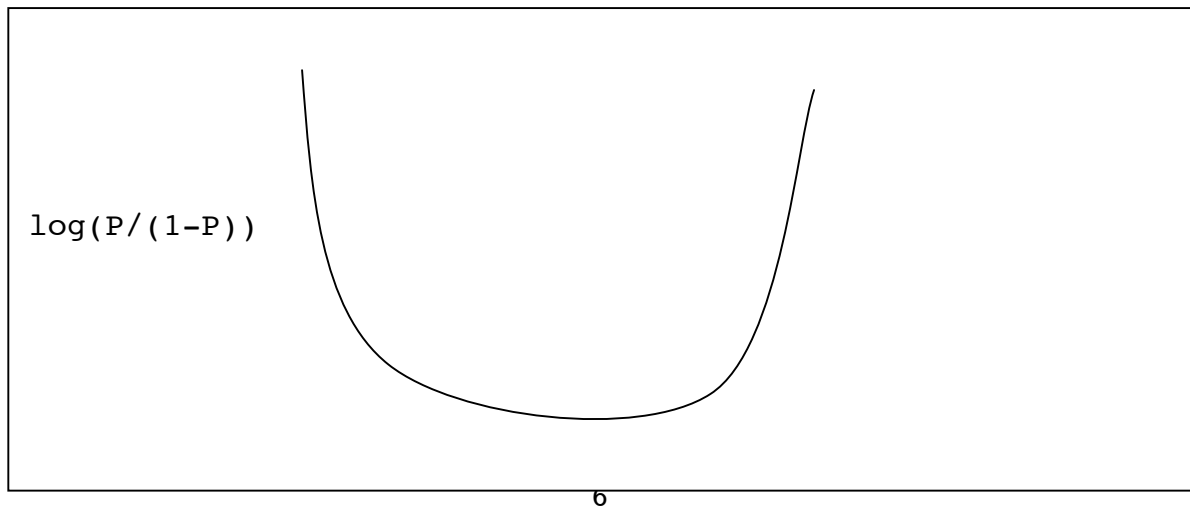
$$\log(P/(1-P)) = B_0^{*} + .0368(AGE)$$

where the value for $B_0^{*}$ depends on the fixed values for the other predictors in the model.

**Figure** Assumed conditional linear relation between age and Log(odds of deaths)



Age

The relationship displayed in this figure refers to an **assumption of linearity**. The slope of the line reflects the increase in the outcome (log odds of dying) per unit increase in the age. For example, the model presented in previous table assumes that the log odds of death increases by 0.0368 for every increase in the year of age. This corresponds to an odds ratio of $e^{.0368} = 1.04$ for each one-year increase in of age.

It is important that the model's assumption properly reflect the true relationship between a continuous predictor and the outcome. If the relationship between a continuous predictor and an outcome is not linear, then the model may need to contain additional terms to reflect its non-linear relationship to the outcome. For example, one might expect that the risk of death might be high for patients with very high blood pressure (hypertension), but also be high for patients with very low blood pressure (hypotension). Thus one might expect a U-shape relationship between blood pressure and the log(odds of death as shown in the following figure:



log(P/(1-P))

Fixing the values of the other predictors, the $MPM_0$ simplifies to the following quadratic equation to describe the conditional association between systolic blood pressure (SBP) and the log(odds of death):

$$\log(P/(1-P)) = B_0^* \; -.0606(SBP) + .000175(SBP^2)$$

where the value for $B_0^*$ depends on the fixed values for the other predictors in the model.

**Relationship Between Stratification and Regression for Controlling Confounding**

The simplest method for controlling a confounder is through stratification according to categories or sub-ranges of the confounder. Strata are sub-groups of patients with common values for the confounder. Since the value for the confounding factor is constant (or nearly constant) within a stratum, all subjects with a stratum should have homogeneous risk of developing the outcome as influenced by the confounder. For example, stratifying by sex will create two strata, one containing only males and the other only females. Even if males may have higher risk for developing the outcome than females, comparing males who received the exposure to males who did not receive the exposure will provide an estimate for the effect of the intervention that is free of confounding by sex. The same is true when the analysis is performed within the female stratum. If the effect of the intervention within males is similar to that within females (no effect modification by gender), then the sex-specific effects can be pooled over strata to provide a single adjusted measure of the effect of the intervention, as demonstrated in a previous series of lecture notes.

This argument can be generalized to confounders with more than two categories to define strata. It can also be applied to continuous confounders by defining strata based on sub-ranges of the confounders. For example, stratification by age often involves strata that are based on decades of age. This approach assumes that the risk of developing the outcome does not vary much by the level of the confounder within each stratum. However, using too wide of a sub-range of the confounder to define a stratum could result in residual confounding within that stratum.

**Control of Multiple Confounders**

Potentially stratification can adjust for the joint confounding by a set of factors through multiple levels of stratification. For example, if age is divided into three age groups (young, middle-aged, and old), then the joint stratification by age and sex will result in six strata (young men, young women, middle-aged men, middle-aged women, old men, and old women). Although simple and straightforward, this method is not practical for adjusting for many confounders. For example, simultaneous stratification by only six binary confounders results in 64 strata, which may be too many for most data sets. Therefore alternative adjustment strategies must be considered. A regression model is the most commonly method for controlling multiple

confounders. This is accomplished by fitting a model containing a term for the intervention/exposure and additional terms for the confounding factors.

If age is treated as a categorical variable (as in the previous paragraph) then controlling for age and sex in a model is analogous to stratification if the model contains separate terms to represent the 6 strata defined by age and sex. However, if age possesses a linear relationship with the log(odds of death) and the effect of age is not modified by sex, then the following simple model may provide a better control of confounding by age and sex than stratification

$$Log(P/(1-P)) = B_0 + B_1(Exposure) + B_2(Age) + B_3(Sex)$$

A regression model that controls for many confounders would be very large and complex. Fixing large models to small or moderate sized data sets is a challenge. One way to resolve this problem is to summarize the confounders into a single score. For example, the following section provides another methods for controlling confounding, not by including terms for the individual confounders in the model, but by summarizing the confounders into summary score, a **propensity score** (*Rosenbaum P, Rubin DB. The central role of the propensity score in observational studies for causal effects. Biometika 1983;70:41-55 and Rosenbaum PR, Rubin DB. Reducing bias in observational studies using subclassification on the propensity score. J Am Stat Assoc 1984;79:516-524.*) If conducted properly, controlling for the propensity score should also control for the individual confounders that define the propensity score.

## Propensity Scores

Typically observational studies (i.e. non-randomized trials) that investigate the effect of clinical interventions (treatments or triage decisions) are characterized by a large number of factors that influence both the choice of the intervention and the outcome. This problem is labeled as **confounding by indication** by epidemiologists, and is demonstrated by the following examples.

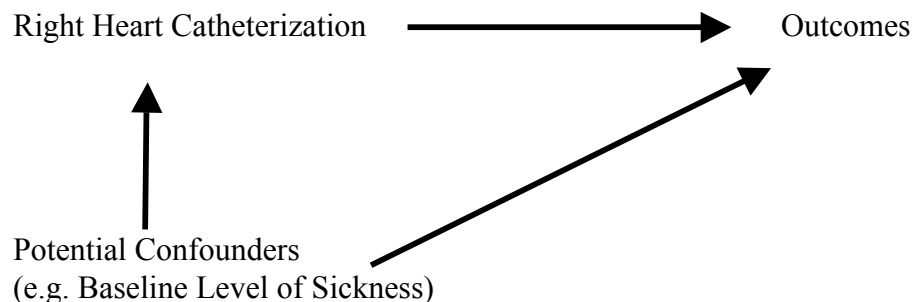**Example # 1: Measuring the Effect of Right Heart Catheterization (SUPPORT)**

Connors et al examined effect of right heart catheterization in the care of critically ill patients (*Connors et al. The effectiveness of right heart catheterization in the initial care of critically ill patients. JAMA. 1996;276:889-897.*). This study examined the survival and health care utilization outcomes for 5735 ICU patients. 2184 of these patients received a right heart catheter (RHC) during the first 24 hours of care in an ICU and another 3551 ICU patients did not receive a RHC. The following table displays the distributions for a sample of patient characteristics.

| Patient Characteristics | Received (RHC) (n=2184) | Did Not Receive RHC (n= 3551) |
|---|---|---|
| Percent over 80 Years of Age | 8% | 14% |
| Mean SBP | 68 | 85 |
| Mean Heart Rate | 119 | 112 |
| Mean Creatinine | 221 | 168 |
| Mean of Apache Score (Measure of Disease Severity) | 61 | 51 |
| Mean Albumin | 29 | 32 |
| Mean of Estimate for 2-Month Survival from Prediction Rule | 56 | 61 |

This table demonstrates that RHC patients differed from the non-RHC patients on a number of factors that can influence outcomes. Therefore, any difference in the outcomes for the two groups of patients might be attributed to the effect of the right heart catheter and/or to the effects of these other factors. This potential for confounding is depicted by the following causal graph (Directed Acyclic Graph, DAG).

**Figure:** Directed Acyclic Graph (DAG) displaying potential for confounding in the study by Connors et al.



Right Heart Catheterization ⟶ Outcomes

Potential Confounders
(e.g. Baseline Level of Sickness)

The following table shows the outcomes (6 month survival, mean total cost, and mean ICU Length of Stay (LOS)) for these patients. Although RHC patients showed worse outcomes, the potential for confounding by the patient characteristics in the previous table raises the question about whether the worse outcomes for the RHC patients reflects the effect of RHC or the type of patient who is given a RHC.

| Outcome | Received (RHC) (n=2184) | Did Not Receive RHC (n= 3551) | P-Value |
|---|---|---|---|
| 6-Month Survival Probability | 53.7% | 46.3% | < 0.001 |
| Mean Total Cost | $131,900 | $74,300 | < 0.001 |
| Mean ICU LOS (days) | 15.5 | 10.3 | < 0.001 |

The following table displays the same patient characteristics for a sample of 1008 RHC and 1008 non-RHC patients who were matched by a propensity score (to be defined later in these notes). Contrary to the previous table for all 2187 RHC patients and 3551 non-RHC patients, the following table shows near identical distributions of these characteristics, suggesting less potential for confounding.
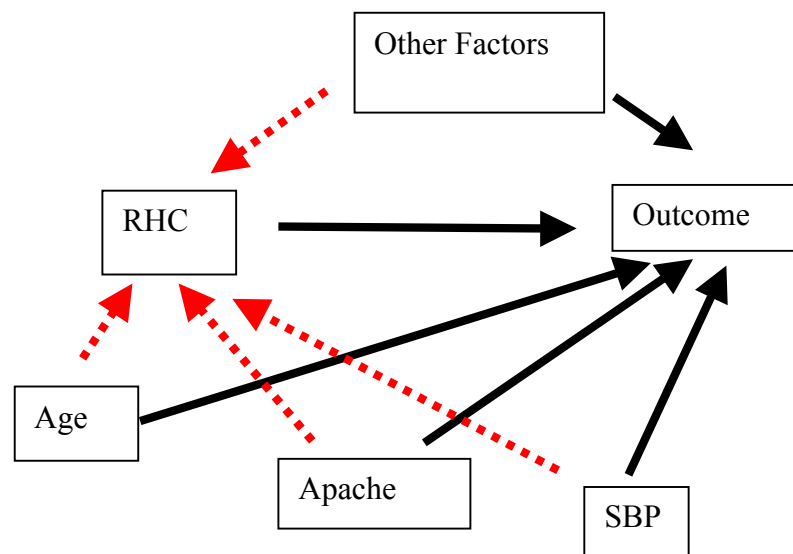
| Patient Characteristics | Received (RHC) (n=2184) | Did Not Receive RHC (n= 3551) |
|---|---|---|
| Mean Age | 60 | 60 |
| Percent Male | 59% | 54% |
| Mean SBP | 71 | 73 |
| Mean Heart Rate | 111 | 111 |
| Mean Creatinine | 203 | 203 |
| Mean of Apache Score (Measure of Disease Severity) | 57 | 57 |
| Mean Albumin | 30 | 30 |
| Mean of Estimate for 2-Month Survival from Prediction Rule | .59 | .58 |

The following table displays the outcomes for the 1008 RHC patients and the 1008 non-RHC patients who were matched by the propensity score. Although the difference in outcomes are attenuated compared to those for all 5735 patients, the outcome for the RHC group remain worse than those for the non-RHC group. These may reflect the true effect of Right Catheterization or possibly confounding by other factors.
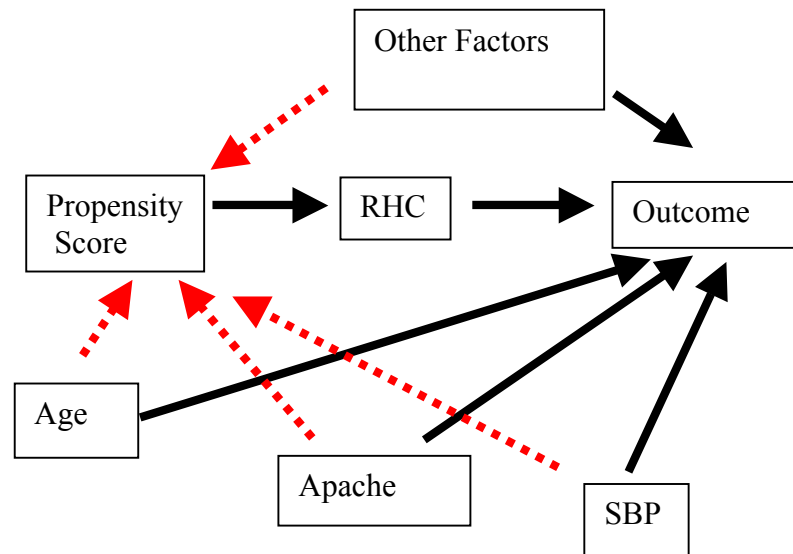
| Outcome | Received (RHC) (n=1008) | Did Not Receive RHC (n= 1008) | P-Value |
|---|---|---|---|
| 6-Month Survival Probability | 51.2% | 46.0% | < 0.001 |
| Mean Total Cost | $49,300 | $35,700 | < 0.001 |
| Mean ICU LOS (days) | 14.8 | 13.0 | < 0.001 |

The motivation for a Propensity Score Analysis is to control for a large number of confounders by combining them into a single summary score. Although the theory of propensity scores was developed over 30 years, their use to control confounding was seldom used until recently.

The Propensity Score is the probability of receiving the exposure as a function of the confounders. The following causal diagram (DAG) displays the anticipated relationship for the factors mentioned in the previous example. The solid arrows connect the suspected confounders and the exposure to the outcome. These arrows reflect the regression coefficients for these predictors in a regression model that predicts the outcome, like the models mentioned earlier in these notes.



On the other hand, the dashed arrows connecting the suspected confounders to the exposure are the basis for the propensity score. The role of the propensity score is to summarize the role of the individual confounders in influencing the treatment decision as shown the following DAG. If the propensity score captures the influence of all of the individual confounders on the treatment decision, then controlling for the propensity score will block the backdoor pathways through the individual factors to the outcome, thereby controlling for any confounding attributed to them.

## Example # 2: Effect of Hypertension Treatment

The following analyses examine the association between hypertension treatment (BPMEDS1) at the 1956 exam and the 24-year risk of death in the FHS teaching data set. The analysis is restricted to 1372 participants with a diagnosis of hypertension at the 1956 exam. The following table shows the crude association between hypertension medication use and death

|  | Died | Survived | Total |
|---|---|---|---|
| BPMEDS1=1 | 91 | 48 | 139 |
| BPMEDS1=0 | 627 | 606 | 1233 |
| Total | 718 | 654 | 1372 |

$$OR_{Crude} = (91/48) / (627/606) = \textbf{1.83}$$

These results suggest a somewhat surprising harmful effect from hypertension medication use. However the following table displays the imbalance of the distribution for 10 potential confounders.

| TABLE | On Hypertension Medication (N=139) | Not on Hypertension Medication (N=1233) |
|---|---|---|
| Male (%) | 30% | 47% |
| Age (Mean) | 56.29 | 53.55 |
| Cholesterol (Mean) | 257.43 | 246.34 |
| SBP (Mean) | 165.08 | 154.26 |
| DBP (Mean) | 96.45 | 93.43 |
| Obese/Overweight (%) | 73% | 72% |
| Smoker (%) | 35% | 42% |
| Diabetes (%) | 6% | 4% |
| Prevalence of CHD (%) | 14% | 7% |
| Prevalence of Stroke(%) | 5% | 1% |

.       The table shows that participants taking hypertension medication have greater average Age, SBP, DBP, Total Cholesterol, Diabetes, History of CHD and Stroke. Participants not on hypertension medication have a higher percentage of Males and Smokers. As mentioned above, the usual method for controlling multiple confounders is by an outcome regression model that contains terms for the individual confounders. For example, the following logistic regression model depicts the risk of death, P, as a function of the exposure (BPMEDS1) and the 10 potential confounding factors

$$\log(P/(1-P)) = -9.4873 + \mathbf{0.3405(BPMED1)} + 0.1117(AGE1) + 0.000239(TOTCHOL1)$$
$$+ 0.0232(SYSBP1) - 0.00894(DIABP1) + 1.2178(MALE)$$
$$+ 0.4763(CURSMOKE1) + 1.6317(DIABETES1) + 1.0511(PRECCHD1)$$
$$+ 0.9972(PREVSTRK1) + 1.3087(UNDERWT) - 0.3166(OVERWT)$$
$$+ 0.0301(OBESE)$$

As described earlier, the coefficient of BPMED1 (0.3405) is the log(Odds Ratio) measuring the association between medication use and death, controlling for the other factors in the model. It follows that the adjusted Odds Ratio is

$$OR_{Logistic} = e^{0.3405} = 1.41$$

An alternative method to control for the potential confounders in this model is by an analysis based on a propensity score, reflecting their relationship of the confounders with the

exposure (BPMEDS1). The following logistic regression model describes the risk of hypertension medication use, PS, as a function of the potential confounders

$$\log(PS/(1-PS)) = -7.0445 + 0.0207(AGE1) + 0.00347(TOTCHOL1)$$
$$+0.0145(SYSBP1) + 0.00592(DIABP1) - 0.4586(MALE)$$
$$- 0.0289(CURSMOKE1) + 0.0697(DIABETES1)$$
$$+ 0.5283(PRECCHD1) + 1.3750(PREVSTRK1)$$
$$+ 0.8119(UNDERWT) + 0.0849(OVERWT) + 0.1080(OBESE)$$

The propensity score is a balancing score. If it captures the relationship between the potential confounders and the exposure, then conditioning on it should eliminate any association between the individual confounders and the exposure. The following table shows the adjusted relationship between the individual confounders and hypertension medication use after adjusting for the propensity score. The first two columns repeat the crude imbalance of the confounders shown in a previously presented table. The last two columns show the balance of the same potential confounders in an analysis that adjusts for the propensity score by we-weighting the data by a function of the propensity score (as demonstrated in the last series of lecture notes). This table demonstrates very similar distributions of the confounders (better balance) after adjusting by the propensity score.

| TABLE | Crude | | Propensity Score Adjusted | |
|---|---|---|---|---|
| | Meds | No Meds | Meds | No Meds |
| Male (%) | 30% | 47% | 43% | 45% |
| Age (Mean) | 56.29 | 53.55 | 53.91 | 53.83 |
| Cholesterol (Mean) | 257.43 | 246.34 | 245.77 | 247.43 |
| SBP (Mean) | 165.08 | 154.26 | 153.86 | 155.37 |
| DBP (Mean) | 96.45 | 93.43 | 92.66 | 93.73 |
| Obese/Overweight (%) | 73% | 72% | 71% | 66% |
| Smoker (%) | 35% | 42% | 41% | 41% |
| Diabetes (%) | 6% | 4% | 4% | 5% |
| Prevalence of CHD (%) | 14% | 7% | 9% | 8% |

| Prevalence of Stroke(%) | 5% | 1% | 2% | 2% |
|---|---|---|---|---|

Confounding can be controlled with a propensity score by the following methods

1. Matching by the propensity score (as in the RHC analysis above)

2. Stratifying by ranges of the propensity score

3. Including the propensity score in an outcome regression model in place of individual confounders

4. Re-weight the data by a function of the propensity score (similar to standardization as described in the previous lecture notes)

The following analysis uses stratification by ranges of the propensity score (second option) to control for confounding in the problem at hand. The following table shows the distribution of the propensity score in seven strata defined by ranges of the propensity score. Not surprising, 33.81% of the participants taking hypertension medication belong to the last four strata, compared to only 12.90% of the participants not on medication.

| Propensity Score | BPMEDS=1 | BPMEDS=0 |
|---|---|---|
| $0 \leq PS \leq 0.05$ | 12 (8.63%) | 234 (18.98%) |
| $0.05 < PS \leq 0.10$ | 43 (30.94%) | 574 (46.55%) |
| $0.10 < PS \leq 0.15$ | 37 (26.62%) | 266 (21.57%) |
| $0.15 < PS \leq 0.20$ | 20 (14.39%) | 82 (6.65%) |
| $0.20 < PS \leq 0.25$ | 9 (6.47%) | 32 (2.60%) |
| $0.25 < PS \leq 0.30$ | 8 (5.76%) | 28 (2.27%) |
| $PS > 0.30$ | 10 (7.19%) | 17 (1.38%) |

Creating seven strata (defined by ranges of the propensity score), measuring the association between hypertension medication use and death within each stratum, and average

these results over the strata using the Mantel-Haenszel formula (presented in the previous series of lecture notes), yields an adjusted Odds Ratio

$$OR_{MH} = 1.37$$

This value for the Odds Ratio from stratifying by the propensity score ($OR_{MH} = 1.37$) is very similar to the adjusted Odds Ratio that was obtained from the outcome logistic regression model ($OR_{Logistic} = 1.41$). The similarity of results is not surprising as both methods are valid options for controlling confounding.

In general, a correctly specified outcome model and an appropriately performed propensity score analysis should result in similar adjusted measures of association to estimate the effect of the exposure. Both methods use regression models, either to prediction the outcome (outcome regression model) or to predict the exposure (propensity score). Using computer simulations, Drake (*Effects of misspecification of the propensity score on estimators of treatment effects. Biometrics, 1993;49:1231-1236*) and Cepeda et al. (*Comparison of logistic regression versus propensity score when the number of events is low and there are multiple confounders. Am J Epidemiol 2003;158:280-287.*) demonstrated that incorrect modeling assumption may have less influence in a propensity score analysis. Cepeda et al. also suggested that a propensity score analysis is superior to an outcome model analysis when the number of outcome events per confounder is small (< 7 events/confounder).

Propensity score analyses are used more frequently that in the past and provide an alternative method for controlling confounding, compared to the commonly used outcome regression model. However, the propensity score analysis (nor the traditional analysis) controls for unmeasured confounders, unless they tend to be highly correlated with those measured confounders present in the model. Even after correctly controlling for multiple confounders by a propensity score analysis (or by an outcome model), the reported adjusted measure of association can still be confounded by other unknown or unmeasured confounders (as might be the case with the adjusted measures of association for the RHC example).