ECON3389 Machine Learning in Economics

Module 2 Classification

Alberto Cappello

Department of Economics, Boston College

Fall 2024

Overview

Agenda:

- Qualitative outcomes and classification problem.
- LPM, logit and probit models.
- Estimation, interpretation and accuracy measures.
- Extensions and other classification algorithms.

Readings:

• ISLR sections 4.1, 4.2, 4.3

- ullet So far our outcome variable Y has always been assumed to be quantitative, e.g. price, quantity, SAT score, etc.
- Qualitative variables (race, gender, geographical region, type of education, etc.) have only been discussed as predictors.
- But what if we want to quantify a relationship where the outcome Y is a qualitative variable?

- ullet So far our outcome variable Y has always been assumed to be quantitative, e.g. price, quantity, SAT score, etc.
- Qualitative variables (race, gender, geographical region, type of education, etc.) have only been discussed as predictors.
- But what if we want to quantify a relationship where the outcome Y is a qualitative variable?
 - A person arrives at an emergency room with symptoms that could be attributed to one of three medical conditions.

- ullet So far our outcome variable Y has always been assumed to be quantitative, e.g. price, quantity, SAT score, etc.
- Qualitative variables (race, gender, geographical region, type of education, etc.) have only been discussed as predictors.
- But what if we want to quantify a relationship where the outcome Y is a qualitative variable?
 - A person arrives at an emergency room with symptoms that could be attributed to one of three medical conditions.
 - Online banking service assesses whether a transaction being performed is fraudulent based on user's IP address, past transaction history, transaction amount, etc.

- ullet So far our outcome variable Y has always been assumed to be quantitative, e.g. price, quantity, SAT score, etc.
- Qualitative variables (race, gender, geographical region, type of education, etc.) have only been discussed as predictors.
- But what if we want to quantify a relationship where the outcome Y is a qualitative variable?
 - A person arrives at an emergency room with symptoms that could be attributed to one of three medical conditions.
 - Online banking service assesses whether a transaction being performed is fraudulent based on user's IP address, past transaction history, transaction amount, etc.
 - A researchers performs an analysis of socio-economic factors that affect whether a student graduates from college or drops out.

Basic Classification Problem

- Consider a qualitative variable Y that for every observation i takes a single value from a finite set of possible unordered values $C = \{y_1, y_2, \dots, y_C\}$.
 - Y = eye color, $C = \{\text{brown}, \text{blue}, \text{green}\}$
 - Y = medical diagnosis, $C = \{\text{stroke}, \text{drug overdose}, \text{epileptic seizure}\}$
 - Y = transaction status, $C = \{\text{fraudulent}, \text{non-fraudulent}\}$

Basic Classification Problem

- Consider a qualitative variable Y that for every observation i takes a single value from a finite set of possible unordered values $C = \{y_1, y_2, \dots, y_C\}$.
 - Y = eye color, $C = \{\text{brown}, \text{blue}, \text{green}\}$
 - Y = medical diagnosis, $C = \{\text{stroke}, \text{drug overdose}, \text{epileptic seizure}\}$
 - Y = transaction status, $C = \{\text{fraudulent}, \text{non-fraudulent}\}$
- Given a feature vector X and a qualitative response Y, the classification task is to build a function C(X) that takes as input the feature vector X and predicts its value for Y, i.e. $C(X) \in \mathcal{C}$.
- In most cases we are interested in estimating the probabilities that Y belongs to a category in \mathcal{C} given X, i.e.

$$Pr(Y = y_c|X) \quad \forall y_c \in C$$



• Suppose for our medical condition classification we code

$$Y = \begin{cases} 1, & \text{if Stroke} \\ 2, & \text{if Drug Overdose} \\ 3, & \text{if Epileptic Seizure} \end{cases}$$

Can we use our standard regression model to predict the medical condition of a patient in the emergency room on the basis of her symptoms?

Suppose for our medical condition classification we code

$$Y = egin{cases} 1, & ext{if Stroke} \ 2, & ext{if Drug Overdose} \ 3, & ext{if Epileptic Seizure} \end{cases}$$

Can we use our standard regression model to predict the medical condition of a patient in the emergency room on the basis of her symptoms?

Does this coding imply an ordering of outcomes?

Suppose for our medical condition classification we code

$$Y = \begin{cases} 1, & \text{if Stroke} \\ 2, & \text{if Drug Overdose} \\ 3, & \text{if Epileptic Seizure} \end{cases}$$

Can we use our standard regression model to predict the medical condition of a patient in the emergency room on the basis of her symptoms?

- Does this coding imply an ordering of outcomes?
- In most cases, it is not possible for us to create a natural ordering in quantitative data

Suppose for our medical condition classification we code

$$Y = \begin{cases} 1, & \text{if Stroke} \\ 2, & \text{if Drug Overdose} \\ 3, & \text{if Epileptic Seizure} \end{cases}$$

Can we use our standard regression model to predict the medical condition of a patient in the emergency room on the basis of her symptoms?

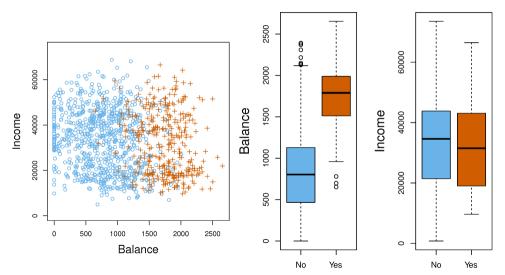
- Does this coding imply an ordering of outcomes?
- In most cases, it is not possible for us to create a natural ordering in quantitative data
- A different coding

$$Y = egin{cases} 1, & ext{if Drug Overdose} \ 2, & ext{if Stroke} \ 3, & ext{if Epileptic Seizure} \end{cases}$$

Will generate a different model and different predictions



Example: Credit Card Defaults



Suppose for our credit card default classification we code

$$Y = egin{cases} 0, & ext{if No} \ 1, & ext{if Yes} \end{cases}$$

Can we use our standard regression model to estimate a regression of Y on X and classify outcome as Yes if $\hat{Y} > 0.5$?

• Suppose for our credit card default classification we code

$$Y = egin{cases} 0, & ext{if No} \ 1, & ext{if Yes} \end{cases}$$

Can we use our standard regression model to estimate a regression of Y on X and classify outcome as Yes if $\hat{Y} > 0.5$?

• Given that Y is binary, we have

$$\mathbb{E}(Y|X) = ?$$



• Suppose for our credit card default classification we code

$$Y = egin{cases} 0, & ext{if No} \ 1, & ext{if Yes} \end{cases}$$

Can we use our standard regression model to estimate a regression of Y on X and classify outcome as Yes if $\hat{Y} > 0.5$?

Given that Y is binary, we have

$$\mathbb{E}(Y|X) = 1 \cdot \mathsf{Pr}(Y = 1|X) + 0 \cdot \mathsf{Pr}(Y = 0|X) = \mathsf{Pr}(Y = 1|X)$$

which means that in this case standard linear regression will estimate the probability of outcome Y = 1, hence the name *linear probability model* or LPM.

• LPM retains all properties of linear regression, but the interpretation of the results is slightly different:

$$\Pr(Y = 1|X) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \ldots + \beta_p X_p + \epsilon$$
$$\beta_j = \frac{\partial \mathbb{E}[\Pr(Y = 1|X)]}{\partial X_j} = \frac{\mathbb{E}[\Delta \Pr(Y = 1|X)]}{\Delta X_j}$$

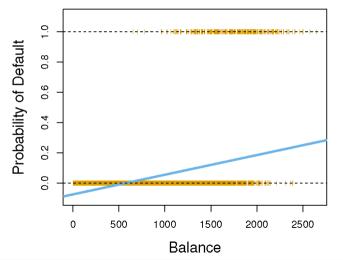
so regression coefficients now capture constant marginal probabilities of outcome Y=1 given a change in X_i .

• LPM retains all properties of linear regression, but the interpretation of the results is slightly different:

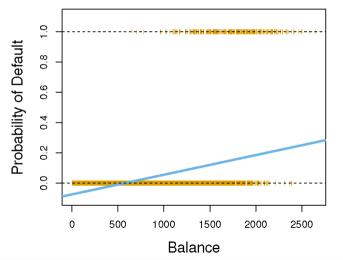
$$\Pr(Y = 1|X) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \ldots + \beta_p X_p + \epsilon$$
$$\beta_j = \frac{\partial \mathbb{E}[\Pr(Y = 1|X)]}{\partial X_j} = \frac{\mathbb{E}[\Delta \Pr(Y = 1|X)]}{\Delta X_j}$$

so regression coefficients now capture constant marginal probabilities of outcome Y=1 given a change in X_i .

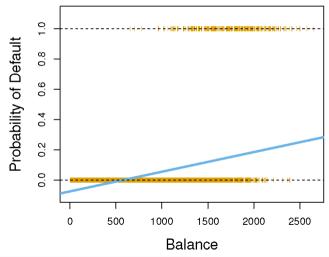
- LPM works very well in binary classification, especially if sample size is moderate to large.
- But it also has some inherent disadvantages, with the most common ones being unreasonable values of $\widehat{\beta}_i$ and predictions for probability outside of [0,1] interval.



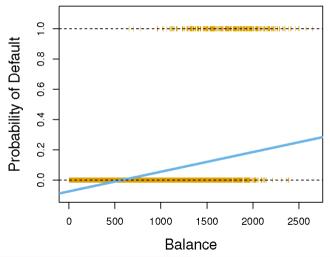
• The orange marks indicate the response Y, either 0 or 1.



- The orange marks indicate the response Y, either 0 or 1.
- Blue line is estimated linear regression, which produces negative predictions for Pr(Y = 1|X) when Balance is less than 500.



- The orange marks indicate the response Y, either 0 or 1.
- Blue line is estimated linear regression, which produces negative predictions for Pr(Y=1|X) when Balance is less than 500.
- One way is to simply ignore the problem and bound any predictions from above and from below.



- The orange marks indicate the response Y, either 0 or 1.
- Blue line is estimated linear regression, which produces negative predictions for Pr(Y=1|X) when Balance is less than 500.
- One way is to simply ignore the problem and bound any predictions from above and from below.
- But can we do better?

Moving away from LPM

• The problem with prediction in LPM stems from the fact that we use linear combination of features X as the estimated probability itself.

Moving away from LPM

- The problem with prediction in LPM stems from the fact that we use linear combination of features X as the estimated probability itself.
- To avoid this problem, we can instead use a one-to-one mapping $F(\cdot)$ from \mathbb{R} to [0,1] interval:

$$\widehat{\Pr}(Y=1|X) = F(\widehat{\beta}_0 + \widehat{\beta}_1 X_1 + \widehat{\beta}_2 X_2 + \ldots + \widehat{\beta}_p X_p)$$

• What could serve as a function $F(\cdot)$? Infinitely many possibilities exist, but because this question was first addressed by statisticians, a natural choice was a *cumulative distribution function* (cdf) from some distribution.

Moving away from LPM

- The problem with prediction in LPM stems from the fact that we use linear combination of features X as the estimated probability itself.
- To avoid this problem, we can instead use a one-to-one mapping $F(\cdot)$ from \mathbb{R} to [0,1] interval:

$$\widehat{\Pr}(Y=1|X) = F(\widehat{\beta}_0 + \widehat{\beta}_1 X_1 + \widehat{\beta}_2 X_2 + \ldots + \widehat{\beta}_p X_p)$$

- What could serve as a function $F(\cdot)$? Infinitely many possibilities exist, but because this question was first addressed by statisticians, a natural choice was a *cumulative distribution function* (cdf) from some distribution.
- For any random variable Z its cdf $F_Z(\cdot)$ is by definition:

$$F_Z(a) = \Pr(Z \leq a)$$

• In classical statistical learning the two most common choice for $F(\cdot)$ are standard normal and logistic cdf

• For simplicity, let's use the matrix notation:

$$\mathbf{X}\boldsymbol{\beta} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \ldots + \beta_p X_p$$

• For simplicity, let's use the matrix notation:

$$\mathbf{X}\boldsymbol{\beta} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \ldots + \beta_p X_p$$

• In probit model $F(\cdot)$ is assumed to be the cdf of a standard normal distribution:

$$F(\boldsymbol{X}eta) = \Phi(\boldsymbol{X}eta) = \int_{-\infty}^{Xeta} rac{1}{\sqrt{2\pi}} \exp^{-rac{z^2}{2}} dz$$

• For simplicity, let's use the matrix notation:

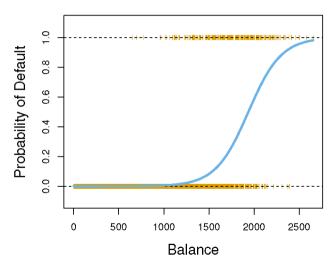
$$\mathbf{X}\boldsymbol{\beta} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \ldots + \beta_p X_p$$

• In probit model $F(\cdot)$ is assumed to be the cdf of a standard normal distribution:

$$F(\boldsymbol{X}\boldsymbol{eta}) = \Phi(\boldsymbol{X}\boldsymbol{eta}) = \int_{-\infty}^{Xeta} \frac{1}{\sqrt{2\pi}} \exp^{-\frac{z^2}{2}} dz$$

• In logit model $F(\cdot)$ is assumed to be the cdf of a logistic distribution:

$$F(\boldsymbol{X}\boldsymbol{eta}) = \Lambda(\boldsymbol{X}\boldsymbol{eta}) = rac{\exp^{\boldsymbol{X}\boldsymbol{eta}}}{1 + \exp^{\boldsymbol{X}\boldsymbol{eta}}}$$



With either normal or logistic cdf as our $F(\cdot)$ function the estimated probability will, by definition, always lie in [0,1] interval.

But if that were the only advantage of logit/probit models, they would not become the de-facto standard econometric model for classification.

• Both logit and probit are special cases of the so-called *generalized linear models* (GLMs). Every GLM consists of three parts: the structural component, the link function and the response distribution.

- Both logit and probit are special cases of the so-called *generalized linear models* (GLMs). Every GLM consists of three parts: the structural component, the link function and the response distribution.
- Structural component is simply the linear combination of our predictors: $X\beta$.

- Both logit and probit are special cases of the so-called *generalized linear models* (GLMs). Every GLM consists of three parts: the structural component, the link function and the response distribution.
- Structural component is simply the linear combination of our predictors: $X\beta$.
- The link function $g(\mu)$ is such that its inverse gives us the (conditional) mean of our outcome Y as a function of the structural component:

$$g(\mu) = oldsymbol{X}oldsymbol{eta}$$
 or $\mathbb{E}(Y|oldsymbol{X}) = \mu = g^{-1}(oldsymbol{X}oldsymbol{eta})$

• The link function is the key to GLMs: since the distribution of the *response variable* Y is non-normal (in our simple example it is binomial), it's what lets us connect the structural component $X\beta$ to the response Y — it 'links' them (hence the name).



- Because our outcome Y is binary, we have $\mathbb{E}(Y|X) = \Pr(Y = 1|X)$, and thus our inverse link function g^{-1} is simply the function that defines conditional probability of Y = 1 given X.
- For probit and logit models the corresponding cumulative distribution functions act as inverse link functions:

$$\mathsf{Probit}: \mathbb{E}(Y|\boldsymbol{X}) = \mu_{Y|X} = \mathsf{Pr}(Y=1|\boldsymbol{X}) = \Phi(\boldsymbol{X}\boldsymbol{\beta})$$

$$\mathsf{Logit}\ : \mathbb{E}(Y|\boldsymbol{X}) = \mu_{Y|X} = \mathsf{Pr}(Y=1|\boldsymbol{X}) = \Lambda(\boldsymbol{X}\boldsymbol{\beta})$$

- Because our outcome Y is binary, we have $\mathbb{E}(Y|X) = \Pr(Y = 1|X)$, and thus our inverse link function g^{-1} is simply the function that defines conditional probability of Y = 1 given X.
- For probit and logit models the corresponding cumulative distribution functions act as inverse link functions:

$$\mathsf{Probit}: \mathbb{E}(Y|\boldsymbol{X}) = \mu_{Y|X} = \mathsf{Pr}(Y=1|\boldsymbol{X}) = \Phi(\boldsymbol{X}\boldsymbol{\beta})$$

$$\mathsf{Logit} \ : \mathbb{E}(Y|\boldsymbol{X}) = \mu_{Y|X} = \mathsf{Pr}(Y=1|\boldsymbol{X}) = \Lambda(\boldsymbol{X}\boldsymbol{\beta})$$

• Note: standard MLR is also a special case of GLM with $g(\mu) = \mu = X\beta$.

- Because our outcome Y is binary, we have $\mathbb{E}(Y|X) = \Pr(Y = 1|X)$, and thus our inverse link function g^{-1} is simply the function that defines conditional probability of Y = 1 given X.
- For probit and logit models the corresponding cumulative distribution functions act as inverse link functions:

Probit :
$$\mathbb{E}(Y|X) = \mu_{Y|X} = \Pr(Y = 1|X) = \Phi(X\beta)$$

Logit : $\mathbb{E}(Y|X) = \mu_{Y|X} = \Pr(Y = 1|X) = \Lambda(X\beta)$

- Note: standard MLR is also a special case of GLM with $g(\mu) = \mu = X\beta$.
- The two key differences of probit/logit models and usual MLR are estimation method and marginal effects calculation/interpretation.

Maximum Likelihood Estimation

• Because we no longer have a direct connection between Y and our structural component $X\beta$, we need to specify our loss function in a different way. Using our link function, we can for every observation i wright down the probability of observing a certain value of Y_i given values of X_i

- Because we no longer have a direct connection between Y and our structural component $X\beta$, we need to specify our loss function in a different way. Using our link function, we can for every observation i wright down the probability of observing a certain value of Y_i given values of X_i
- For example, for a logit model we have:

$$\mathsf{Pr}(Y = Y_i | \boldsymbol{X}_i) = \left(\frac{\mathsf{exp}^{X_i\beta}}{1 + \mathsf{exp}^{X_i\beta}}\right)^{Y_i} \left(1 - \frac{\mathsf{exp}^{X_i\beta}}{1 + \mathsf{exp}^{X_i\beta}}\right)^{1 - Y_i}$$

- Because we no longer have a direct connection between Y and our structural component $X\beta$, we need to specify our loss function in a different way. Using our link function, we can for every observation i wright down the probability of observing a certain value of Y_i given values of X_i
- For example, for a logit model we have:

$$\Pr(Y = Y_i | \boldsymbol{X}_i) = \left(\frac{\exp^{X_i\beta}}{1 + \exp^{X_i\beta}}\right)^{Y_i} \left(1 - \frac{\exp^{X_i\beta}}{1 + \exp^{X_i\beta}}\right)^{1 - Y_i}$$

• With the default assumption of i.i.d. observations we can wright down the joint probability or *likelihood function* of seeing our sample:

$$\ell(oldsymbol{eta}) = \prod_{i=1}^n \mathsf{Pr}(Y = Y_i | oldsymbol{X}_i)$$



• $Maximum\ likelihood\ estimation\ (ML)$ is a method that chooses parameters β so as to minimize the loss function in form of the negative of the likelihood function:

$$\widehat{oldsymbol{eta}}_{ extit{ML}} = \mathop{\mathsf{argmin}}_{oldsymbol{eta}} - \ell(oldsymbol{eta})$$

• $Maximum\ likelihood\ estimation\ (ML)$ is a method that chooses parameters eta so as to minimize the loss function in form of the negative of the likelihood function:

$$\widehat{oldsymbol{eta}}_{ extit{ML}} = \mathop{\mathsf{argmin}}_{oldsymbol{eta}} - \ell(oldsymbol{eta})$$

- Under some general conditions $\widehat{\beta}_{ML}$ is efficient, consistent and asymptotically normal, just like $\widehat{\beta}_{OLS}$. In fact, one can show that standard OLS is a special case of ML if the error term ϵ in MLR is exactly normal.
- ullet But unlike OLS, ML is a more general estimation procedure and allows one to recover structural parameters such as eta in models that are far more flexible than standard MLR.

Marginal Effects in Probit/Logit

 The other key difference of logit/probit models from LPM is the fact that margial effects are now calculated and interpreted in a different way:

Probit :
$$\frac{\partial p(\mathbf{X})}{\partial X_j} = \frac{\partial \Phi(\mathbf{X}\beta)}{\partial X_j} \beta_j \neq \beta_j$$
 Logit : $\frac{\partial p(\mathbf{X})}{\partial X_j} = \frac{\partial \Lambda(\mathbf{X}\beta)}{\partial X_j} \beta_j \neq \beta_j$

where $p(\boldsymbol{X}) = \Pr(Y = 1 | \boldsymbol{X})$ for simplicity

Marginal Effects in Probit/Logit

• The other key difference of logit/probit models from LPM is the fact that margial effects are now calculated and interpreted in a different way:

Probit :
$$\frac{\partial p(\mathbf{X})}{\partial X_j} = \frac{\partial \Phi(\mathbf{X}\beta)}{\partial X_j} \beta_j \neq \beta_j$$
 Logit : $\frac{\partial p(\mathbf{X})}{\partial X_j} = \frac{\partial \Lambda(\mathbf{X}\beta)}{\partial X_j} \beta_j \neq \beta_j$

where $p(X) = \Pr(Y = 1|X)$ for simplicity

- The marginal effects now depend on values of all variables in X, so we need to either estimate the marginal effects at a specific value of all our predictors (typically means or medians) or calculate their average over all values of X in our sample.
- ullet Additionally, structural parameters eta no longer have any direct interpretation on their own, with the exception of a few special cases.

- While LPM can use the standard R^2 as a well-defined goodness-of-fit measure, it is not an option for logit/probit models due to the different loss function.
- In standard MLR an single R^2 value of 0.95 is an evidence of an excellent fit, but in classification problems we are often more interested in *class-specific* performance, especially in areas such as medicine or biology.

- While LPM can use the standard R^2 as a well-defined goodness-of-fit measure, it is not an option for logit/probit models due to the different loss function.
- In standard MLR an single R^2 value of 0.95 is an evidence of an excellent fit, but in classification problems we are often more interested in *class-specific* performance, especially in areas such as medicine or biology.
- Consider the case where Y=1 means a positive test result. Then our model's predictions fall into one of the 4 possible cases:

[Y = 0	Y = 1
	$\widehat{Y} = 0$		
	$\widehat{Y}=1$		

- While LPM can use the standard R^2 as a well-defined goodness-of-fit measure, it is not an option for logit/probit models due to the different loss function.
- In standard MLR an single R^2 value of 0.95 is an evidence of an excellent fit, but in classification problems we are often more interested in *class-specific* performance, especially in areas such as medicine or biology.
- Consider the case where Y=1 means a positive test result. Then our model's predictions fall into one of the 4 possible cases:

	Y = 0	Y = 1
$\widehat{Y}=0$	True Negative	
$\widehat{Y}=1$		

- While LPM can use the standard R^2 as a well-defined goodness-of-fit measure, it is not an option for logit/probit models due to the different loss function.
- In standard MLR an single R^2 value of 0.95 is an evidence of an excellent fit, but in classification problems we are often more interested in *class-specific* performance, especially in areas such as medicine or biology.
- Consider the case where Y=1 means a positive test result. Then our model's predictions fall into one of the 4 possible cases:

	Y = 0	Y = 1		
$\widehat{Y} = 0$	True Negative	False Negative (Type II Error)		
$\widehat{Y}=1$				

- While LPM can use the standard R^2 as a well-defined goodness-of-fit measure, it is not an option for logit/probit models due to the different loss function.
- In standard MLR an single R^2 value of 0.95 is an evidence of an excellent fit, but in classification problems we are often more interested in *class-specific* performance, especially in areas such as medicine or biology.
- Consider the case where Y=1 means a positive test result. Then our model's predictions fall into one of the 4 possible cases:

	Y = 0	Y = 1
$\widehat{Y} = 0$	True Negative	False Negative (Type II Error)
$\widehat{Y}=1$	False Positive (Type I Error)	

- While LPM can use the standard R^2 as a well-defined goodness-of-fit measure, it is not an option for logit/probit models due to the different loss function.
- In standard MLR an single R^2 value of 0.95 is an evidence of an excellent fit, but in classification problems we are often more interested in *class-specific* performance, especially in areas such as medicine or biology.
- ullet Consider the case where Y=1 means a positive test result. Then our model's predictions fall into one of the 4 possible cases:

	Y = 0	Y=1	
$\widehat{Y} = 0$	True Negative	False Negative (Type II Error)	
$\widehat{Y}=1$	False Positive (Type I Error)	True Positive	

- While LPM can use the standard R^2 as a well-defined goodness-of-fit measure, it is not an option for logit/probit models due to the different loss function.
- In standard MLR an single R^2 value of 0.95 is an evidence of an excellent fit, but in classification problems we are often more interested in *class-specific* performance, especially in areas such as medicine or biology.
- Consider the case where Y=1 means a positive test result. Then our model's predictions fall into one of the 4 possible cases:

	Y = 0	Y = 1	
$\widehat{Y} = 0$	True Negative	False Negative (Type II Error)	
$\widehat{Y}=1$	False Positive (Type I Error)	True Positive	

- While LPM can use the standard R^2 as a well-defined goodness-of-fit measure, it is not an option for logit/probit models due to the different loss function.
- In standard MLR an single R^2 value of 0.95 is an evidence of an excellent fit, but in classification problems we are often more interested in *class-specific* performance, especially in areas such as medicine or biology.
- ullet Consider the case where Y=1 means a positive test result. Then our model's predictions fall into one of the 4 possible cases:

	Y = 0	Y = 1		
$\widehat{Y} = 0$ True Negative		False Negative (Type II Error)		
$\widehat{Y}=1$	False Positive (Type I Error)	True Positive		

• For a COVID-19 test or cancer screening, we care more FN then about FP.

- While LPM can use the standard R^2 as a well-defined goodness-of-fit measure, it is not an option for logit/probit models due to the different loss function.
- In standard MLR an single R^2 value of 0.95 is an evidence of an excellent fit, but in classification problems we are often more interested in *class-specific* performance, especially in areas such as medicine or biology.
- Consider the case where Y=1 means a positive test result. Then our model's predictions fall into one of the 4 possible cases:

	Y = 0	Y = 1	
$\widehat{Y} = 0$ True Negative		False Negative (Type II Error)	
$\widehat{Y}=1$	False Positive (Type I Error)	True Positive	

- For a COVID-19 test or cancer screening, we care more FN then about FP.
- For city administration FPs in traffic cameras and speeding tickets are more important.

- While LPM can use the standard R^2 as a well-defined goodness-of-fit measure, it is not an option for logit/probit models due to the different loss function.
- In standard MLR an single R^2 value of 0.95 is an evidence of an excellent fit, but in classification problems we are often more interested in *class-specific* performance, especially in areas such as medicine or biology.
- Consider the case where Y=1 means a positive test result. Then our model's predictions fall into one of the 4 possible cases:

	Y = 0	Y=1	
$\widehat{Y} = 0$ True Negative		False Negative (Type II Error)	
$\widehat{Y}=1$	False Positive (Type I Error)	True Positive	

- For a COVID-19 test or cancer screening, we care more FN then about FP.
- For city administration FPs in traffic cameras and speeding tickets are more important.
- In judicial system both FP and FN are equally important.

		True default status		
		No	Yes	Total
	No	9644	252	9896
Predicted default status	Yes	23	81	104
	Total	9667	333	10000

		True default status		
		No	Yes	Total
	No	9644	252	9896
Predicted default status	Yes	23	81	104
	Total	9667	333	10000

- If we simply look at pure prediction precision, then:
 - Our total error rate is (23 + 252)/10000 = 2.75%, which seems low enough.

		True default status		
		No	Yes	Total
	No	9644	252	9896
Predicted default status	Yes	23	81	104
	Total	9667	333	10000

- If we simply look at pure prediction precision, then:
 - Our total error rate is (23 + 252)/10000 = 2.75%, which seems low enough.
 - Out of 104 predicted defaults 81 ended up being classified correctly, which means only 23/9667 = 0.24% of all non-defaults were classified incorrectly.

		True default status		
		No	Yes	Total
Predicted default status	No	9644	252	9896
	Yes	23	81	104
	Total	9667	333	10000

- If we simply look at pure prediction precision, then:
 - Our total error rate is (23 + 252)/10000 = 2.75%, which seems low enough.
 - Out of 104 predicted defaults 81 ended up being classified correctly, which means only 23/9667 = 0.24% of all non-defaults were classified incorrectly.
 - However, out of 333 true defaults we managed to miss 252/333 = 75.67%, which could be an unacceptably high error rate for this class.

	Y = 0	Y = 1
$\widehat{Y} = 0$	True Negative Rate (TNR) or specificity:	
	TNR = TN/N = 9644/9667 = 99.76%	
$\widehat{Y}=1$		

	Y = 0	Y = 1
$\widehat{Y} = 0$	True Negative Rate (TNR) or specificity:	False Negative Rate (FNR):
	TNR = TN/N = 9644/9667 = 99.76%	FNR = FN/P = 252/333 = 75.67%
$\widehat{Y}=1$		

	Y = 0	Y = 1
$\widehat{Y}=0$	True Negative Rate (TNR) or specificity:	False Negative Rate (FNR):
	TNR = TN/N = 9644/9667 = 99.76%	FNR = FN/P = 252/333 = 75.67%
$\widehat{Y}=1$	False Positive Rate (FPR):	
	FPR = FP/N = 23/9667 = 0.24%	

	Y = 0	Y = 1
$\widehat{Y} = 0$	True Negative Rate (TNR) or specificity:	False Negative Rate (FNR):
	TNR = TN/N = 9644/9667 = 99.76%	FNR = FN/P = 252/333 = 75.67%
$\widehat{Y}=1$	False Positive Rate (FPR):	True Positive Rate (TPR) or sensitivity:
	FPR = FP/N = 23/9667 = 0.24%	TPR = TP/P = 81/333 = 24.33%

	Y = 0	Y = 1
$\widehat{Y} = 0$	True Negative Rate (TNR) or specificity:	False Negative Rate (FNR):
	TNR = TN/N = 9644/9667 = 99.76%	FNR = FN/P = 252/333 = 75.67%
$\widehat{Y}=1$	False Positive Rate (FPR):	True Positive Rate (TPR) or sensitivity:
	FPR = FP/N = 23/9667 = 0.24%	TPR = TP/P = 81/333 = 24.33%

• This why in classification problems it is important to evaluate class-specific precision via the following four measures:

	Y = 0	Y = 1
$\widehat{Y}=0$	True Negative Rate (TNR) or specificity:	False Negative Rate (FNR):
	TNR = TN/N = 9644/9667 = 99.76%	FNR = FN/P = 252/333 = 75.67%
$\widehat{Y}=1$	False Positive Rate (FPR):	True Positive Rate (TPR) or sensitivity:
	FPR = FP/N = 23/9667 = 0.24%	TPR = TP/P = 81/333 = 24.33%

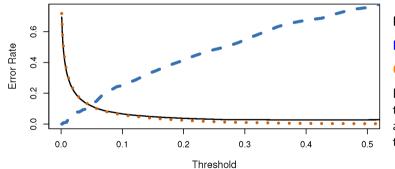
• As one can easily see, all four measures are related to each other. In particular, the following two identities must always hold:

$$\mathsf{TNR} + \mathsf{FPR} = 100\%$$
 and $\mathsf{TPR} + \mathsf{FNR} = 100\%$



The table on the previous slide was constructed using the rule $\hat{Y} = 1$ if $\Pr(Y = 1 | X) > 0.5$, because 0.5 is the most common probability threshold used for classification predictions. However, the values of all 4 goodness-of-fit metrics will change if we change this threshold.

The table on the previous slide was constructed using the rule $\widehat{Y} = 1$ if $\widehat{\Pr}(Y = 1 | X) > 0.5$, because 0.5 is the most common probability threshold used for classification predictions. However, the values of all 4 goodness-of-fit metrics will change if we change this threshold.



Black line: total error rate

Blue dashes: FNR

Orange dots: FPR

Based on this chart, we might want to set our threshold to 0.05 to achieve better error rate composition.

While logit and probit models usually deliver very similar estimation results (especially on large datasets), modern statistical learning overwhelmingly prefers to use logistic regression. Why?

While logit and probit models usually deliver very similar estimation results (especially on large datasets), modern statistical learning overwhelmingly prefers to use logistic regression. Why?

• Coefficient interpretation. In Economics we are interested in calculating and interpreting marginal effects, but in logit model one can also interpret the actual values of $\widehat{\beta}_j$ themselves. This is because in logit model coefficient $\widehat{\beta}_j$ shows how the log of odds ratio changes with changes in X_j :

$$ln\left(\frac{p(\mathbf{X})}{1-p(\mathbf{X})}\right) = \mathbf{X}\boldsymbol{\beta} \quad \Rightarrow \quad \beta_j = \frac{\partial ln\left(\frac{p(X)}{1-p(X)}\right)}{\partial X_j}$$

While logit and probit models usually deliver very similar estimation results (especially on large datasets), modern statistical learning overwhelmingly prefers to use logistic regression. Why?

• Coefficient interpretation. In Economics we are interested in calculating and interpreting marginal effects, but in logit model one can also interpret the actual values of $\widehat{\beta}_j$ themselves. This is because in logit model coefficient $\widehat{\beta}_j$ shows how the log of odds ratio changes with changes in X_j :

$$ln\left(\frac{p(\boldsymbol{X})}{1-p(\boldsymbol{X})}\right) = \boldsymbol{X}\boldsymbol{\beta} \quad \Rightarrow \quad \beta_j = \frac{\partial ln\left(\frac{p(\boldsymbol{X})}{1-p(\boldsymbol{X})}\right)}{\partial X_j}$$

• Random utility models. Suppose consumer is choosing between two alternatives based on utility that is a function of observable product attributes X and a random utility shock ϵ . Then if ϵ follows Type I EV distribution, consumer's choice probabilities will take logit form (McFadden, D. (1973)).

While logit and probit models usually deliver very similar estimation results (especially on large datasets), modern statistical learning overwhelmingly prefers to use logistic regression. Why?

• Coefficient interpretation. In Economics we are interested in calculating and interpreting marginal effects, but in logit model one can also interpret the actual values of $\widehat{\beta}_j$ themselves. This is because in logit model coefficient $\widehat{\beta}_j$ shows how the log of odds ratio changes with changes in X_j :

$$ln\left(\frac{p(\mathbf{X})}{1-p(\mathbf{X})}\right) = \mathbf{X}\boldsymbol{\beta} \quad \Rightarrow \quad \beta_j = \frac{\partial ln\left(\frac{p(\mathbf{X})}{1-p(\mathbf{X})}\right)}{\partial X_j}$$

- Random utility models. Suppose consumer is choosing between two alternatives based on utility that is a function of observable product attributes X and a random utility shock ϵ . Then if ϵ follows Type I EV distribution, consumer's choice probabilities will take logit form (McFadden, D. (1973)).
- Generalized choice models and information theory (Matejka, F. and McKay, A. (2015).

Multinomial and Ordered Outcomes

• In binary outcome case all three options (LPM, logit, probit) are applicable and most times yield similar results, especially in large samples. But once we move beyond basic binary case, LPM becomes structurally infeasible.

Multinomial and Ordered Outcomes

- In binary outcome case all three options (LPM, logit, probit) are applicable and most times yield similar results, especially in large samples. But once we move beyond basic binary case, LPM becomes structurally infeasible.
- Two hospitals use the following coding for incoming ER patients:

$$Y = \begin{cases} 1, & \text{if drug overdose} \\ 2, & \text{if stroke} \\ 3, & \text{if epileptic seizure} \end{cases} \quad \text{and} \quad Y = \begin{cases} 1, & \text{if stroke} \\ 5, & \text{if epileptic seizure} \\ 6, & \text{if drug overdose} \end{cases}$$

• This is an example of *multinomial unordered* classification. The fundamental difference with binary case is that here different coding will yield different results in LPM, but not in logit/probit.

Multinomial and Ordered Outcomes

- In binary outcome case all three options (LPM, logit, probit) are applicable and most times yield similar results, especially in large samples. But once we move beyond basic binary case, LPM becomes structurally infeasible.
- Two hospitals use the following coding for incoming ER patients:

$$Y = \begin{cases} 1, & \text{if drug overdose} \\ 2, & \text{if stroke} \\ 3, & \text{if epileptic seizure} \end{cases}$$
 and $Y = \begin{cases} 1, & \text{if stroke} \\ 5, & \text{if epileptic seizure} \\ 6, & \text{if drug overdose} \end{cases}$

- This is an example of *multinomial unordered* classification. The fundamental difference with binary case is that here different coding will yield different results in LPM, but not in logit/probit.
- Same thing happens with ordered outcomes, e.g. "strongly disagree, disagree, uncertain, agree, strongly agree", where we cannot impose a standardized numerical difference between outcomes (as opposed to something like number of kids in the family as the outcome).

Other SL and ML Classification Algorithms

- While research in econometrics over the past 50 years has developed a very wide range of discrete choice models, it wasn't just economics and casual inference that has been driving the development of statistical learning in classification problems.
- Areas such as genetics, biostatistics, pharmaceutics and others have always had a need for statistical models that could precisely classify certain outcomes.

Other SL and ML Classification Algorithms

- While research in econometrics over the past 50 years has developed a very wide range of discrete choice models, it wasn't just economics and casual inference that has been driving the development of statistical learning in classification problems.
- Areas such as genetics, biostatistics, pharmaceutics and others have always had a need for statistical models that could precisely classify certain outcomes.
- In recent decade huge advances in new variations of previously less used methods such as neural networks have been achieved due to ever-increasing demand for classification and prediction in modern data-dominated areas such as image and voice recognition.
- We will not be covering those methods in details (at least not till later in the course), because all of them are nearly completely irrelevat for causal inference. But because they have some other advantages, they deserve an honorary mentioning.

Discriminant analysis

• Instead of directly modeling Pr(Y|X), model the distribution of X in each of the C classes separately, and then use Bayes theorem to flip things around and obtain Pr(Y|X):

$$\Pr(Y = j | X = x) = \frac{\Pr(Y = j) \Pr(X = x | Y = j)}{\sum_{j=1}^{C} \Pr(Y = j) \Pr(X = x | Y = j)} = \frac{\pi_j f_j(x)}{\sum_{j=1}^{C} \pi_j f_j(x)}$$

where $f_j(x) = \Pr(X = x | Y = j)$ is the *density* for X in class j and $\pi_j = \Pr(Y = j)$ is the *prior* probability for class j.

Discriminant analysis

• Instead of directly modeling Pr(Y|X), model the distribution of X in each of the C classes separately, and then use *Bayes theorem* to flip things around and obtain Pr(Y|X):

$$\Pr(Y = j | X = x) = \frac{\Pr(Y = j) \Pr(X = x | Y = j)}{\sum_{j=1}^{C} \Pr(Y = j) \Pr(X = x | Y = j)} = \frac{\pi_j f_j(x)}{\sum_{j=1}^{C} \pi_j f_j(x)}$$

where $f_j(x) = \Pr(X = x | Y = j)$ is the *density* for X in class j and $\pi_j = \Pr(Y = j)$ is the *prior* probability for class j.

• Both $f_j(x)$ and π_j are estimated from the data, with the most common choice being normal density for $f_j(x)$.

Discriminant analysis

• Instead of directly modeling Pr(Y|X), model the distribution of X in each of the C classes separately, and then use Bayes theorem to flip things around and obtain Pr(Y|X):

$$\Pr(Y = j | X = x) = \frac{\Pr(Y = j) \Pr(X = x | Y = j)}{\sum_{j=1}^{C} \Pr(Y = j) \Pr(X = x | Y = j)} = \frac{\pi_j f_j(x)}{\sum_{j=1}^{C} \pi_j f_j(x)}$$

where $f_j(x) = \Pr(X = x | Y = j)$ is the *density* for X in class j and $\pi_j = \Pr(Y = j)$ is the *prior* probability for class j.

- Both $f_j(x)$ and π_j are estimated from the data, with the most common choice being normal density for $f_j(x)$.
- ullet DA works especially well in small samples with well-separated classes and X variables that have approximately normal distributions.

K-nearest Neighbors

• For any given test observation x_0 and a positive integer K, first identify K points in the training data that are closest to x_0 , denoted as \mathcal{N}_0 . Then the conditional probability for class j is calculated as

$$\Pr(Y = j | X = x_0) = \frac{1}{K} \sum_{i \in \mathcal{N}_0} \mathbb{I}(y_i = j)$$

which is simply a fraction of points in \mathcal{N}_0 with response values equal to j.

K-nearest Neighbors

• For any given test observation x_0 and a positive integer K, first identify K points in the training data that are closest to x_0 , denoted as \mathcal{N}_0 . Then the conditional probability for class j is calculated as

$$\Pr(Y = j | X = x_0) = \frac{1}{K} \sum_{i \in \mathcal{N}_0} \mathbb{I}(y_i = j)$$

which is simply a fraction of points in \mathcal{N}_0 with response values equal to j.

- KNN is an example of non-parametric supervised learning algorithm and as such places almost no restrictions on the nature of the data.
- It quickly loses its potency when the number of features in X grows above 4-5 (too many points in high-dimensional space could be equally close to x_0). Thus it is often paired with other methods aimed at feature extraction and dimensionality reduction, such as *principal component analysis* (PCA).

Decision trees

• Create a sequence of binary splits that partitions (stratifies, segments) the feature space X (i.e. the set of possible values for X_1, X_2, \ldots, X_p) into J distinct and non-overlapping regions R_1, R_2, \ldots, R_J . Then predict that each observation belongs to the *most commonly occurring class* of training observations in the region to which it belongs.

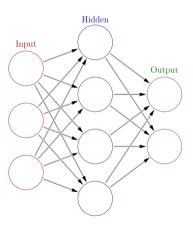
Decision trees

- Create a sequence of binary splits that partitions (stratifies, segments) the feature space X (i.e. the set of possible values for X_1, X_2, \ldots, X_p) into J distinct and non-overlapping regions R_1, R_2, \ldots, R_J . Then predict that each observation belongs to the *most commonly occurring class* of training observations in the region to which it belongs.
- Decision trees are very easy to explain and interpret, especially when they are small enough to be displayed graphically.
- Some believe that decision trees are a more natural way to model human decision making, as opposed to regressions and other methods.

Decision trees

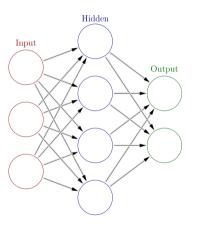
- Create a sequence of binary splits that partitions (stratifies, segments) the feature space X (i.e. the set of possible values for X_1, X_2, \ldots, X_p) into J distinct and non-overlapping regions R_1, R_2, \ldots, R_J . Then predict that each observation belongs to the *most commonly occurring class* of training observations in the region to which it belongs.
- Decision trees are very easy to explain and interpret, especially when they are small enough to be displayed graphically.
- Some believe that decision trees are a more natural way to model human decision making, as opposed to regressions and other methods.
- In their basic forms decision trees tend to have inferior levels of prediction accuracy compared to even basic LPM, but by employing modifications that allow aggregation and randomization of many trees into ensembles (forests) one can improve their accuracy drastically, although at the cost of interpretability.

Neural Networks



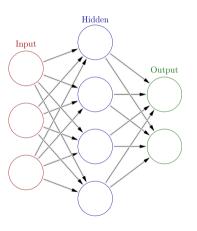
Artificial neural networks (ANNs) is a learning/computing system
that vaguely resembles neural systems in animal brains. ANNs consist of multiple layers (input layer, output layer and one or more
hidden layers) of nodes called *neurons*, each capable of receiving and
transmitting signals via connections to other neurons.

Neural Networks



- Artificial neural networks (ANNs) is a learning/computing system
 that vaguely resembles neural systems in animal brains. ANNs consist of multiple layers (input layer, output layer and one or more
 hidden layers) of nodes called *neurons*, each capable of receiving and
 transmitting signals via connections to other neurons.
- Each connection transferring the output of a neuron to the input of another neuron is assigned a weight. The propagation function computes the input to a neuron from the outputs of its predecessor neurons and their connections as a weighted sum.

Neural Networks



- Artificial neural networks (ANNs) is a learning/computing system
 that vaguely resembles neural systems in animal brains. ANNs consist of multiple layers (input layer, output layer and one or more
 hidden layers) of nodes called *neurons*, each capable of receiving and
 transmitting signals via connections to other neurons.
- Each connection transferring the output of a neuron to the input of another neuron is assigned a weight. The propagation function computes the input to a neuron from the outputs of its predecessor neurons and their connections as a weighted sum.
- ANN learns by adjusting its weighted associations according to a learning rule, using the error between the inputs and predicted outputs.