

# **INTRO TO DATA SCIENCE**

## **LECTURE 7: PROBABILITY AND NAIVE BAYES CLASSIFICATION**

**Paul Burkard**

**07/16/2015**

## **LAST TIME:**

- CLASSIFICATION**
- K-NEAREST NEIGHBORS CLASSIFICATION**

## **QUESTIONS?**

**I INTRO TO PROBABILITY**

**II. BAYESIAN INFERENCE**

**III. NAIVE BAYES CLASSIFICATION**

**HANDS-ON: NAIVE BAYES CLASSIFICATION**

**IV. LOGISTIC REGRESSION**

# **I. INTRO TO PROBABILITY**

*Q: What is a **probability**?*

*A: A number between 0 and 1 that characterizes the likelihood that some event will occur.*

*The probability of event  $A$  is denoted  $P(A)$ .*

*Q: What is the set of all possible events called?*

*A: This set is called the **sample space**  $\Omega$ . Event  $A$  is a member of the sample space, as is every other event.*

*The probability of the sample space  $P(\Omega)$  is 1.*

*Q: What is a probability distribution?*

*A: A function that assigns probability to each event in the sample space.*

*A distribution can be **discrete** or **continuous***

*Ex: Discrete – Uniform distribution*

$$X \sim \{1, \dots, N\} \longrightarrow P(X = x) = 1/N$$

*Continuous – Normal distribution –  $N(\mu, \sigma^2)$*

$$X \sim N(0, 1) \longrightarrow P(X = x) = 0$$

***Discrete Probability Distributions:***

- These are ***probability mass functions***
- Each value  $P(X=x)$  represents the probability of observing a given value  $x$  for variable  $X$

$$P(\Omega) = \sum_{X=x} P(X = x) = 1$$

$$P(\Omega) = \sum_{X=x} P(X = x) = 1$$



***Continuous Probability Distributions:***

- These are **probability density functions (PDF)**
- Each value  $P(X=x)$  represents the **relative probability** of observing a given value  $x$  for variable  $X$

$$P(x_0 < x < x_1) = \int_{x_0}^{x_1} P(x) dx$$

$$P(\Omega) = \int_{-\infty}^{+\infty} P(x) dx = 1$$

*Q: What is a random variable?*

*A: Essentially, a measurable whose possible values have a probability distribution*

*Values of these are the “Events” for which we’re looking to measure the probabilities*

*Q: What is expected value?*

*A: It is the average value of a random variable*

*For discrete distributions*

$$E(X) = \sum x * p(x)$$

*For continuous distributions*

$$E(X) = \int (x * p(x)) dx$$

*Q: Consider two events  $A$  &  $B$ . How can we characterize the intersection of these events?*

*A: With the joint probability of  $A$  and  $B$ , written  $P(A, B)$ .*

*Q: Suppose event  $B$  has occurred. What quantity represents the probability of  $A$  **given** this information about  $B$ ?*

*A: The intersection of  $A$  &  $B$  divided by region  $B$ .*

*This is called the **conditional probability** of  $A$  given  $B$ , written  $P(A|B) = P(A \cap B) / P(B)$ .*

*Q: What does it mean for two events to be independent?*

*A: Information about one does not affect the probability of the other.*

*This can be written as  $P(A|B) = P(A)$ .*

*Using the definition of the conditional probability, we can also write:*

$$P(A|B) = P(A \cap B) / P(B) = P(A) \rightarrow P(A \cap B) = P(A) * P(B)$$

# **II. BAYESIAN INFERENCE**

*This result is called Bayes' theorem. Here it is again:*

$$P(A|B) = P(B|A) * P(A) / P(B)$$

*Some facts:*

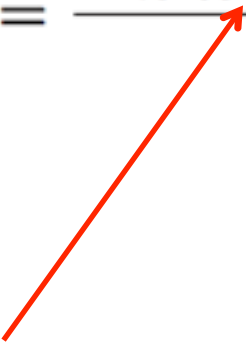
- This is a simple algebraic relationship using elementary definitions.*
- It's a very powerful computational tool.*



*Each term in this relationship has a name, and each plays a distinct role in any probability calculation (including ours). Here's how it might look in the context of classification.*

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$

*This term is the **likelihood function**. It represents the joint probability of observing features  $\{x_i\}$  given that that record belongs to class  $C$ .*

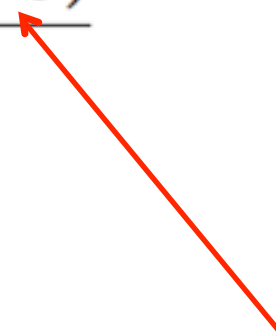
$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$


*This term is the **likelihood function**. It represents the joint probability of observing features  $\{x_i\}$  given that that record belongs to class  $C$ .*

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$

*We can observe the value of the likelihood function from the training data.*

*This term is the **prior probability** of  $C$ . It represents the probability of a record belonging to class  $C$  before the data is taken into account.*

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$


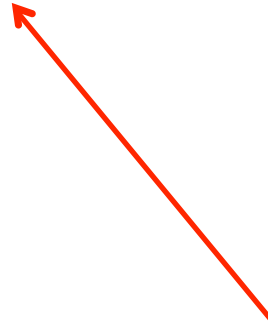
*This term is the **prior probability** of  $C$ . It represents the probability of a record belonging to class  $C$  before the data is taken into account.*

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$

*The value of the prior is also observed from the data.*

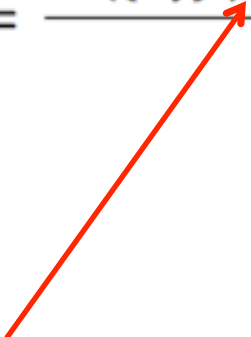
*This term is the **normalization constant**. It doesn't depend on  $C$ , and is generally ignored until the end of the computation.*

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$



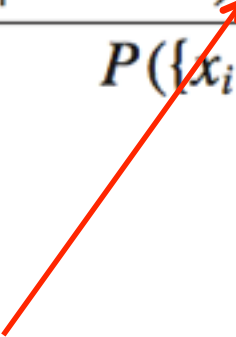
*Maximum likelihood estimator (MLE):*

*What parameters **maximize** the likelihood function?*

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$


*Maximum a posteriori estimate (MAP):*

*What parameters **maximize** the likelihood function **AND** prior?*

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$




# **III. NAIVE BAYES CLASSIFICATION**

	<i><b>continuous</b></i>	<i><b>categorical</b></i>
<i><b>supervised</b></i>	<i>regression</i>	<i>classification</i>
<i><b>unsupervised</b></i>	<i>dim reduction</i>	<i>clustering</i>

*This term is the **posterior probability** of  $C$ . It represents the probability of a record belonging to class  $C$  after the data is taken into account.*

$$\boxed{P(\text{class } C \mid \{x_i\})} = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$

*The goal of any Bayesian computation is to find (“learn”) the posterior distribution of a particular variable given new **evidence**.*

Problem:

We observe the following coin flips:

HTHH

What is  $P(X = \text{Heads})$  ?

Problem:

We observe the following coin flips:

HTHH

What is  $P(X = \text{Heads})$  ?  $3/4$ , Why?

Problem:

We observe the following coin flips:

HTHHTHT

What is  $P(X = \text{Heads})$  ?

Problem:

We observe the following coin flips:

HTHHTHT

What is  $P(X = \text{Heads})$  ?  $4/7$ , Why?

We observe the following coin flips:

HTHHTHT

*Maximum likelihood estimator (MLE):*

*What parameters **maximize** the likelihood function?*

Let  $P(X = \text{Heads}) = q$ , and write Bayes Theorem

$$P(q \mid \text{observations}) = P(\text{observations} \mid q) * P(q) / \text{constant}$$



*Maximum likelihood estimator (MLE):*

*What parameters **maximize** the likelihood function?*

Let  $P(X = \text{Heads}) = q$ , and write Bayes Theorem

$$P(q \mid \text{observations}) = P(\text{observations} \mid q) * \\ P(q) / \text{constant}$$

$$P(\text{observations} \mid q) = ?$$

$$P(q) = ?$$

Binomial Distribution:

$$\Pr(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$$

$$\begin{aligned} P(4 \text{ heads, } 3 \text{ tails} \mid q) &= P(X = 4, n = 7) \\ &= (7 \text{ choose } 4) * q^4 * (1-q)^3 \end{aligned}$$

Binomial Distribution:

$$\Pr(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$$

$$\begin{aligned} P(4 \text{ heads, } 3 \text{ tails} \mid q) &= P(X = 4, n = 7) \\ &= (7 \text{ choose } 4) * q^4 * (1-q)^3 \end{aligned}$$

Optimize w.r.t.  $q \longrightarrow$  **MLE:**  $q = 4/7$

A prior distribution is known as **conjugate prior** if its from the same family as the posterior for a certain likelihood function

For the binomial distribution, the conjugate prior is the **Beta distribution**

$$\begin{aligned} &= \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1 - x)^{\beta-1} \\ &= \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1 - x)^{\beta-1} \end{aligned}$$

The **MAP estimate** is the value that maximizes both the likelihood function and prior – the product of the two.

In the coin flip setting is the value that optimizes

$$P(4H, 3T \mid q) * P(q)$$

The **MAP estimate** is the value that maximizes both the likelihood function and prior – the product of the two.

In the coin flip setting with a Beta distribution it's the value that optimizes:

$$\begin{aligned} &P(4H, 3T \mid q) * P(q) \\ &= \binom{7}{4} q^4 (1-q)^3 q^{(\alpha-1)} (1-q)^{(\beta-1)} \end{aligned}$$

Why do you care?

Why do you care?

Many problems are binary and are estimated using counts...



Why do you care?

Many problems are binary and are estimated using counts...

Ex. 1:

Sample 100 people and ask if they support a politician?

Why do you care?

Many problems are binary and are estimated using counts...

Ex. 1:

Sample 100 people and ask if they support a politician?

23 say Yes – Is the correct prediction 23/100?

What's the prior?

Ex. 2:

Need to choose between multiple categories to present (for ads, products, news).

Ex. 2:

Need to choose between multiple categories to present (for ads, products, news).

You can compute response % for each category

Ex. 2:

Need to choose between multiple categories to present (for ads, products, news).

You can compute response % for each category

But each should have a unique prior – **unique psuedo counts**

*Suppose we have a dataset with features  $x_1, \dots, x_n$  and a class label  $C$ .  
What can we say about classification using Bayes' theorem?*

*Suppose we have a dataset with features  $x_1, \dots, x_n$  and a class label  $C$ . What can we say about classification using Bayes' theorem?*

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$

*Bayes' theorem can help us to determine the probability of a record belonging to a class, given the data we observe.*

*The idea of Bayesian inference, then, is to **update** our beliefs about the distribution of  $C$  using the data (“evidence”) at our disposal.*

$$P(\text{class } C \mid \{x_i\}) = \frac{P(\{x_i\} \mid \text{class } C) \cdot P(\text{class } C)}{P(\{x_i\})}$$

*Then we can use the posterior for prediction.*



*Q: What piece of the puzzle we've seen so far looks like it could be intractably difficult in practice?*

*Remember the likelihood function?*

$$P(\{x_i\} | C) = P(\{x_1, x_2, \dots, x_n\} | C)$$

*Remember the likelihood function?*

$$P(\{x_i\} | C) = P(\{x_1, x_2, \dots, x_n\} | C)$$

*Observing this exactly would require us to have enough data for every possible combination of features to make a reasonable estimate.*

*Q: What piece of the puzzle we've seen so far looks like it could intractably difficult in practice?*

*A: Estimating the full likelihood function.*

*Q: So what can we do about it?*

*Q: So what can we do about it?*

*A: Make a simplifying assumption. In particular, we assume that the features  $x_i$  are conditionally independent from each other:*

*Q: So what can we do about it?*

*A: Make a simplifying assumption. In particular, we assume that the features  $x_i$  are conditionally independent from each other:*

$$P(\{x_i\} | C) = P(x_1, x_2, \dots, x_n | C) \approx P(x_1 | C) * P(x_2 | C) * \dots * P(x_n | C)$$

*Q: So what can we do about it?*

*A: Make a simplifying assumption. In particular, we assume that the features  $x_i$  are conditionally independent from each other:*

$$P(\{x_i\} | C) = P(x_1, x_2, \dots, x_n | C) \approx P(x_1 | C) * P(x_2 | C) * \dots * P(x_n | C)$$

*This “**naïve**” **assumption** simplifies the likelihood function to make it tractable.*



$$P(\{x_i\} | C) = P(x_1, x_2, \dots, x_n | C) \approx P(x_1 | C) * P(x_2 | C) * \dots * P(x_n | C)$$

*Q: Given that we can compute this value, what do we do with it?*

$$P(\{x_i\} | C) = P(x_1, x_2, \dots, x_n | C) \approx P(x_1 | C) * P(x_2 | C) * \dots * P(x_n | C)$$

*Q: Given that we can compute this value, what do we do with it?*

*A: In our training phase, we ‘learn’ the probability of seeing our training examples under each class.*

$$P(\{x_i\} | C) = P(x_1, x_2, \dots, x_n | C) \approx P(x_1 | C) * P(x_2 | C) * \dots * P(x_n | C)$$

*Q: Given that we can compute this value, what do we do with it?*

*A: In our training phase, we 'learn' the probability of seeing our training examples under each class.*

*Then we use Bayes Theorem to compute  $P(\text{class} | \text{inputs})$*

*Example: Text Classification*

***Does this news article talk about politics?***

*Training Set: Collection of New Articles*

*Example: Text Classification*

***Does this news article talk about politics?***

*Training Set: Collection of New Articles*

*Article 1: The computer contractor who exposed....*

*Article 2: The parents of a missing U.S. journalist in Syria...*

*Q: What are my features?*

*Q: What are my features?*

*A: The text in the documents.*

*Q: What are my features?*

*A: The text in the documents.*

*Q: How to I represent them?*



*Q: What are my features?*

*A: The text in the documents.*

*Q: How do I represent them?*

*A: Binary occurrence? Word counts?*

*the, computer, contractor, exposed, parents, missing, Syria, U.S.*

<i>1</i>	<i>1</i>	<i>1</i>	<i>1</i>	<i>0</i>	<i>0</i>	<i>0</i>	<i>0</i>
<i>1</i>	<i>0</i>	<i>0</i>	<i>0</i>	<i>1</i>	<i>1</i>	<i>1</i>	<i>1</i>

*the, computer, contractor, exposed, parents, missing, Syria, U.S.*

<i>1</i>	<i>1</i>	<i>1</i>	<i>1</i>	<i>0</i>	<i>0</i>	<i>0</i>	<i>0</i>
<i>1</i>	<i>0</i>	<i>0</i>	<i>0</i>	<i>1</i>	<i>1</i>	<i>1</i>	<i>1</i>

*We can make some alterations*

*1) Drop stop words (commonly occurring words that don't have meaning)*

*the, computer, contractor, exposed, parents, missing, Syria, U.S., **POL***

*1      1      1      1      0      0      0      0*

*1      0      0      0      1      1      1      1*

*Our goal is to compute*

*$P ( POL = T \mid \text{words in the text} )$*

*We need to **learn**  $P( \text{word} \mid POL )$*

*i.e.  $P ( Syria \mid POL )$*

*the, computer, contractor, exposed, parents, missing, Syria, U.S., **POL***

*1      1      1      1      0      0      0      0*

*1      0      0      0      1      1      1      1*

*Once we've learned  $P(\text{computer} \mid \text{POL})$ ,  $P(\text{U.S.} \mid \text{POL})$  etc. on our training set, we want to label our test set*

*the, computer, contractor, exposed, parents, missing, Syria, U.S., **POL***

*1      1      1      1      0      0      0      0*

*1      0      0      0      1      1      1      1*

*The predicted label,  $POL = \text{True}$  or*

*$POL = \text{False}$  is the one that maximizes our posterior.*

*the, computer, contractor, exposed, parents, missing, Syria, U.S., **POL***

*1          1          1          1          0          0          0          0*

*1          0          0          0          1          1          1          1*

*Compute probability in each class:*

$$P ( POL = T \mid \{x\} ) = c * P ( \{x\} \mid POL = T ) * P(POL=T)$$

$$P ( POL = F \mid \{x\} ) = c * P ( \{x\} \mid POL = F ) * P(POL=F)$$

*the, computer, contractor, exposed, parents, missing, Syria, U.S., **POL***

*1          1          1          1          0          0          0          0*

*1          0          0          0          1          1          1          1*

*Article 2: The parents of a missing U.S. journalist in Syria...*

$$P ( POL = T \mid \{x\} ) = P ( \{x\} \mid POL = T ) * P(POL=T)$$

$$= P(Syria \mid POL=T) * P(journalist \mid POL=T) * P(parents \mid POL=T) ... * P( POL=T)$$



# **IV. LOGISTIC REGRESSION**

*Q: What is logistic regression?*

***Q: What is logistic regression?***

***A: A generalization of the linear regression model to classification problems.***

*In linear regression, we used a set of covariates to predict the value of a (continuous) outcome variable.*

*In linear regression, we used a set of covariates to predict the value of a (continuous) outcome variable.*

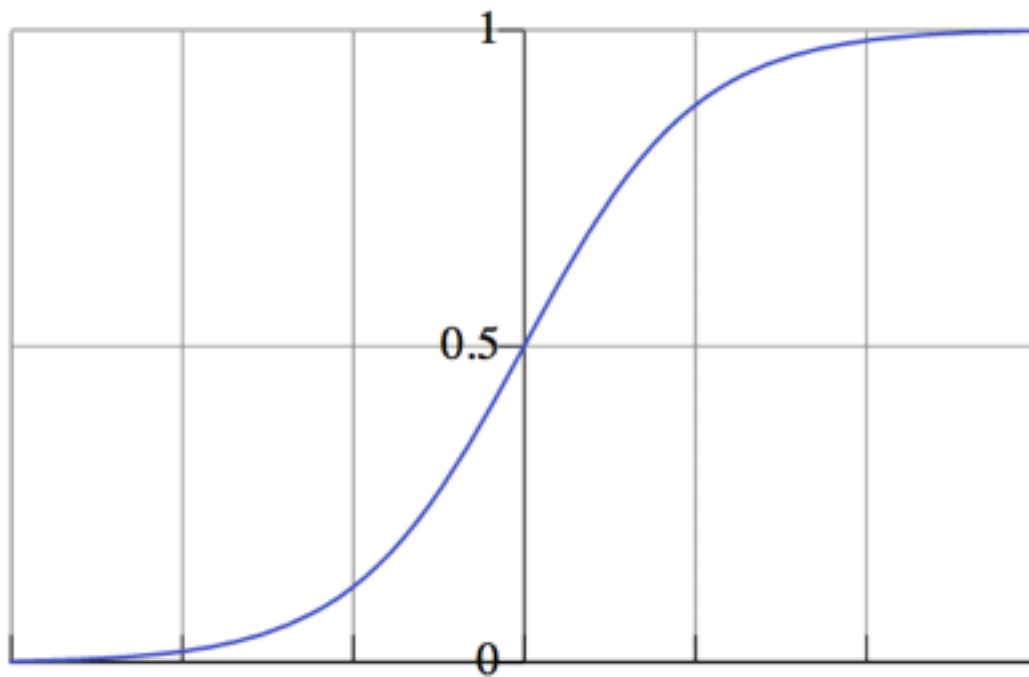
*In logistic regression, we use a set of covariates to predict probabilities of (binary) class membership.*

*In linear regression, we used a set of covariates to predict the value of a (continuous) outcome variable.*

*In logistic regression, we use a set of covariates to predict probabilities of (binary) class membership.*

*These probabilities are then mapped to class labels, thus solving the classification problem.*

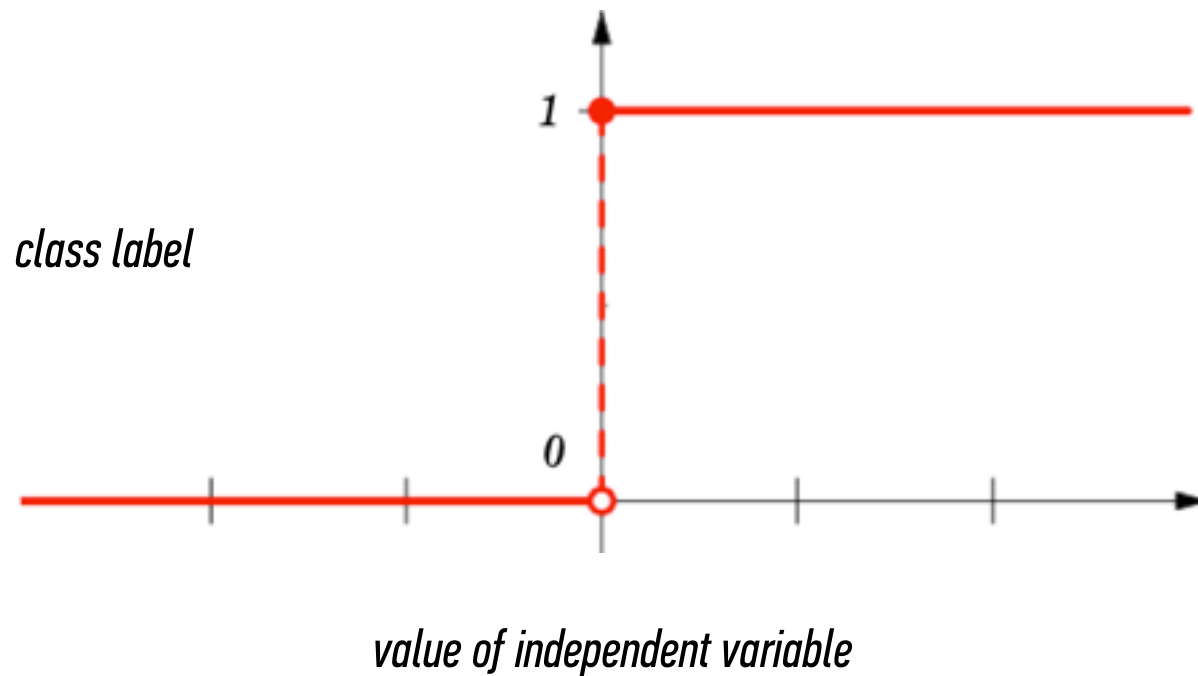
*probability of  
belonging to  
class*



*value of independent variable*

## NOTE

Probability predictions look like this.



## NOTE

Probabilities are “snapped” to class labels (eg by thresholding at 50%).



*The logistic regression model is an extension of the linear regression model, with a couple of important differences.*

*The logistic regression model is an extension of the linear regression model, with a couple of important differences.*

*The main difference is in the outcome variable.*

*The key variable in any regression problem is the **response type** of the outcome variable  $y$  given the value of the covariate  $x$ :*

$$E(y|x)$$

*The key variable in any regression problem is the **conditional mean** of the outcome variable  $y$  given the value of the covariate  $x$ :*

*In linear regression, we assume  $E(y|x)$  conditional mean is a linear function taking values in  $(-\infty, +\infty)$ :*

$$E(y|x) = \alpha + \beta x$$

*In logistic regression, we've seen that the conditional mean of the outcome variable takes values only in the unit interval  $[0, 1]$ .*

*In logistic regression, we've seen that the conditional mean of the outcome variable takes values only in the unit interval  $[0, 1]$ .*

*The first step in extending the linear regression model to logistic regression is to map the outcome variable  $E(y|x)$  into the unit interval.*

*In logistic regression, we've seen that the conditional mean of the outcome variable takes values only in the unit interval  $[0, 1]$ .*

*The first step in extending the linear regression model to logistic regression is to map the outcome variable  $E(y \mid x)$  into the unit interval.*

*Q: How do we do this?*

*A: By using a transformation called the **logistic function**:*

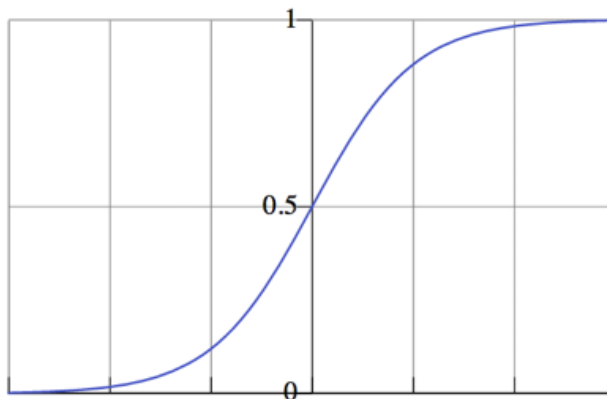
$$E(y|x) = \pi(x) = \frac{e^{\alpha + \beta x}}{1 + e^{\alpha + \beta x}}$$



*A: By using a transformation called the **logistic function**:*

*We've already seen what this looks like.*

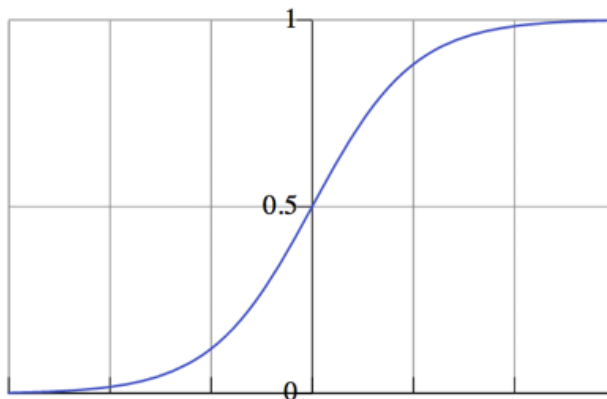
$$E(y|x) = \pi(x) = \frac{e^{\alpha + \beta x}}{1 + e^{\alpha + \beta x}}$$



*A: By using a transformation called the **logistic function**:*

$$E(y|x) = \pi(x) = \frac{e^{\alpha + \beta x}}{1 + e^{\alpha + \beta x}}$$

*We've already seen how to use the logistic function.*



## NOTE

For any value of  $x$ ,  $y$  is in the interval  $[0, 1]$

This is a nonlinear transformation!

*The **logit function** is an important transformation of the logistic function. Notice that it returns the linear model!*

$$g(x) = \ln\left(\frac{\pi(x)}{1-\pi(x)}\right) = \alpha + \beta x$$

*The **logit function** is an important transformation of the logistic function. Notice that it returns the linear model!*

$$g(x) = \ln\left(\frac{\pi(x)}{1-\pi(x)}\right) = \alpha + \beta x$$

*The logit function is also called the **log-odds function**.*

*The **logit function** is an important transformation of the logistic function. Notice that it returns the linear model!*

$$g(x) = \ln\left(\frac{\pi(x)}{1-\pi(x)}\right) = \alpha + \beta x$$

*The logit function is also called the **log-odds function**.*

### NOTE

This name hints at its usefulness in interpreting our results.

We will see why shortly.

- *Classification Problems*
- *When we need an estimate of class likelihood (“probabilistic classifier”)*
- *Many attributes*