

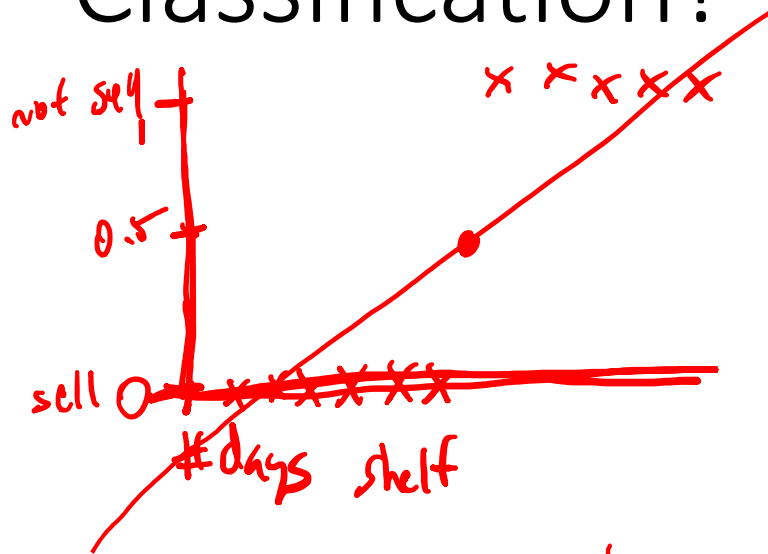
# Logistic Regression

# Logistic Regression

independent      dependent  
 $y = mx + b$

- Logistic Regression is a terrible name!
  - It's not regression at all!
  - It's classification
- But as you'll see, how we do it is extremely similar to *linear* regression

# Linear Regression for Classification?

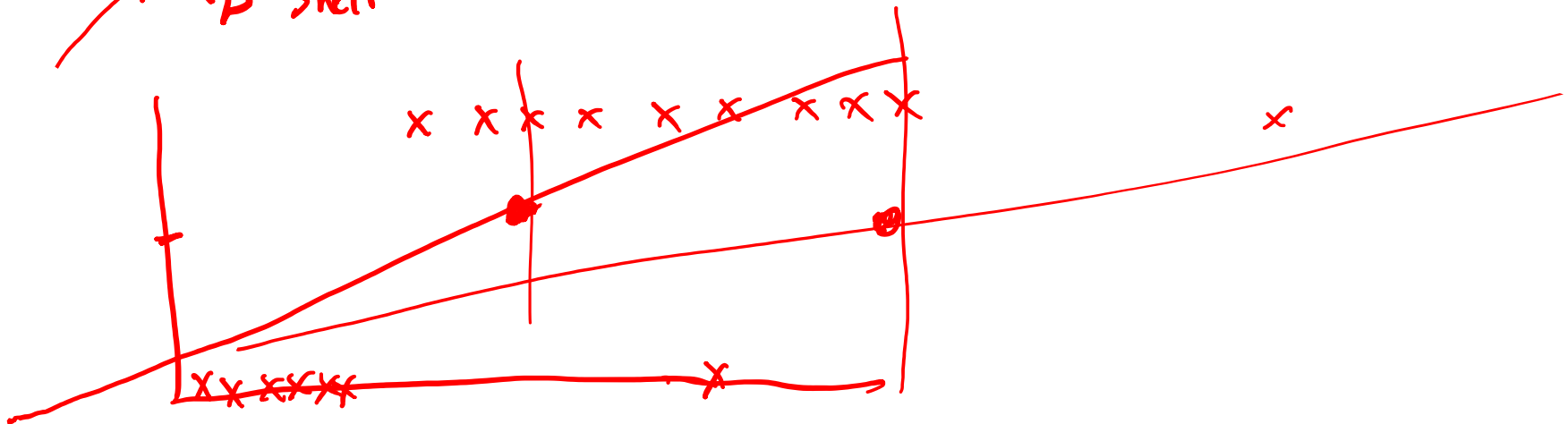


$$y \in \{0, 1\}$$

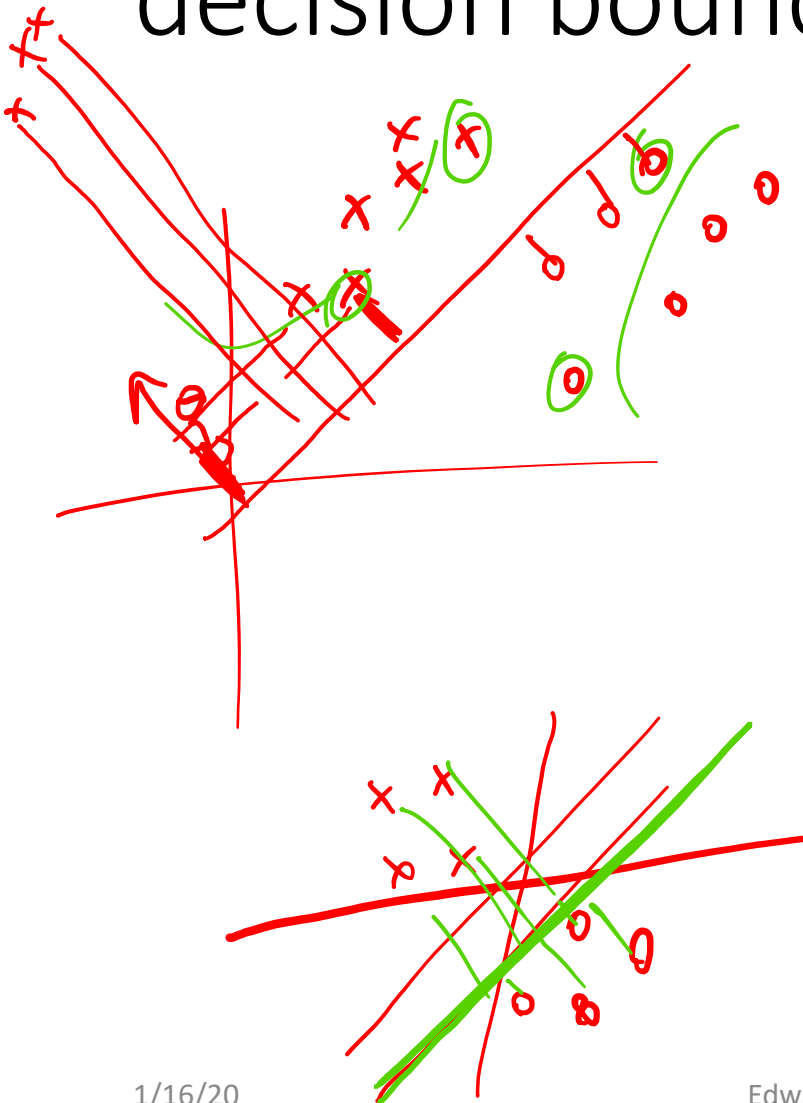
sell no

$$\hat{y} \geq 0.5 \rightarrow 1$$

$$\hat{y} < 0.5 \rightarrow 0$$



# Logistic Regression – Line as decision boundary



$$y \in \{0, 1\}$$

$$ax + by - c = 0$$

$$\theta^T x = 0$$

$$\underline{\theta \cdot x_i} = \text{distance to boundary}$$

$$\underline{\theta \cdot p_i = 1000 \text{ units?}}$$

$$\underline{\theta \cdot x_i} \Rightarrow \text{Prob Cat } \underline{[0-1]}$$

# Logistic Regression

- With logistic regression we assume binary classification and want to provide a probability for the positive class:

$$0 \leq \underline{P(y = 1)} \leq 1$$

- Recall from *linear* regression we computed

$$g(x, \theta) = \underline{x\theta}$$

- We can alter this for use in computing  $P(y = 1)$  as:

$$P(y = 1) = g(x, \theta) = \frac{1}{1 + e^{-x\theta}}$$

sigmoid  
logistic  
squashing

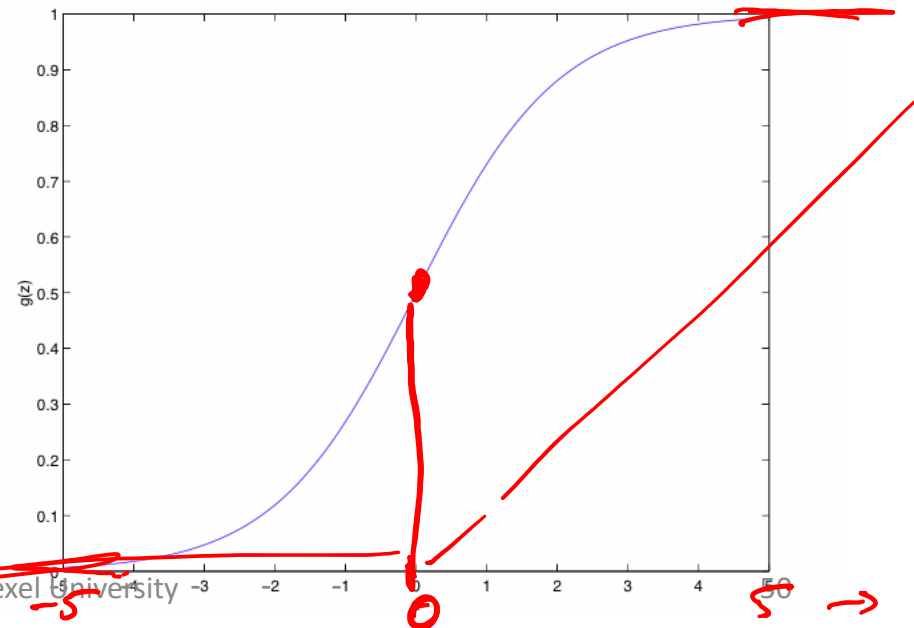
# Logistic Regression

$$P(y = 1|x) = g(x, \theta) = \frac{1}{1 + e^{-x\theta}}$$

$\theta x \rightarrow 1000$   
 $\rightarrow .9999$

$\theta x = -1000$   
 $\rightarrow 0$

- This function, Let  $g(z) = \frac{1}{1+e^{-z}}$  is called the *sigmoid* or *logistic* function
  - Tends to 0 as  $z$  decreases
  - Tends to 1 as  $z$  increases
- This has the nice characteristic in that it's differentiable
  - Why might that be important?!



# Logistic Regression



- If we consider

- $P(y = 1 | x, \theta) = \frac{1}{1 + e^{-x\theta}} = g(x, \theta)$

- Then we can compute the probability of being from the negative class as:

- $P(y = 0 | x, \theta) = 1 - g(x, \theta)$

- Ultimately we want to find the parameters  $\theta$  to minimize the classification error

- Or conversely, to find the parameters  $\theta$  to maximize the correct class likelihood

$$P(y | x, \theta) = g(x, \theta)^y (1 - g(x, \theta))^{(1-y)}$$

$$y = 1 \Rightarrow g(x, \theta)$$

$$y = 0 \Rightarrow (1 - g(x, \theta))$$

# Fit Parameters Based on Maximum Likelihood

- MLE*  
*Maximum Likelihood Estimation*
- Given a supervised observation  $(x, y)$ , we can compute the **likelihood** that we are correct as  $\ell(y|x, \theta) = \underline{(g(x, \theta))^y (1 - g(x, \theta))^{(1-y)}}$

Doing this for the entire dataset, since the observations are conditionally independent of one another we have:

$$\ell(Y|X, \theta) = \prod_{t=1}^N \ell(Y_t, |X_t, \theta) = \prod_{t=1}^N (g(X_t, \theta)^{Y_t} (1 - g(X_t, \theta))^{(1-Y_t)})$$

*iid*  $\rightarrow$  *independently identically distributed*



$$\log(ab) = \log a + \log b$$

$$\log(a^b) = b \log a$$

# Log Likelihood

- So what do we do with this likelihood  $\ell(Y|X, \theta)$ ?
  - We want to maximize it!
  - Or minimize  $-\ell(Y|X, \theta)$
- So we're going to want to take the derivative
- But taking the derivative of a product of a lot of things involves a very long expansion
- Let's instead first take the *log* of this
  - Doing so will result in a sum which is easier to take the derivative of.
  - So now we want to maximize the log likelihood

$$\log(a^b) = b \log a$$

$$\log(ab) = \log a + \log b$$

# Log Likelihood Rules

$$\prod_{i=1}^N g(x_i, \theta)^{y_i} \cdot (1 - g(x_i, \theta))^{(1-y_i)}$$

$$\sum_{i=1}^N y_i \ln g(x_i, \theta) + (1 - y_i) \ln (1 - g(x_i, \theta))$$

# Log Likelihood

- From the properties of logarithms
  - $\log_b(mn) = \log_b(m) + \log_b(n)$
  - $\log_b(m^n) = n \cdot \log_b(m)$
- Returning to our likelihood for a single observation,  $\ell(y, |x, \theta)$ , we get
  - $\ell(y|x, \theta) = \ln \left( g(x, \theta)^y (1 - g(x, \theta))^{(1-y)} \right)$
  - $= \ln(g(x, \theta)^y) + \ln \left( (1 - g(x, \theta))^{1-y} \right)$
  - $= y \ln(g(x, \theta)) + (1 - y) \ln(1 - g(x, \theta))$

# Log Likelihood

- Since we're taking the log of product of this for each instance we get a sum!

$$\ell(Y|X, \theta) = \log P(Y|X, \theta) = \sum_{t=1}^N Y_t \ln(g(X_t, \theta)) + (1 - Y_t) \ln(1 - g(X_t, \theta))$$

- Ideally we'd like to take the derivative of this with respect to  $\theta$ , set it equal to zero, and solve for  $\theta$  to find the maxima
  - The ~~closed form approach~~
  - But this doesn't exist ☹️
- So what's our other approach?
  - Do partial derivatives on the parameters and use gradient descent! (actually in this case gradient ascent, since we're ~~trying to~~ maximize)
  - First do this for a single observation and single parameter
    - Then vectorize!

# Log Likelihood Derivation

$$\left( \frac{y}{g(x, \theta)} - \frac{1-y}{1-g(x, \theta)} \right) \cdot x_j \cdot \cancel{g(x, \theta) - (1-g(x, \theta))}$$

$$\theta = \theta - \frac{1}{N} X^T (y - g(x, \theta))$$

$x_{\theta}$

$$\frac{(1-g(x, \theta))y}{(1-g(x, \theta))(g(x, \theta))} - \frac{g(x, \theta)(1-y)}{g(x, \theta)(1-g(x, \theta))}$$

$$\frac{\partial J}{\partial \theta_j} = (y - g(x, \theta)) \cdot x_j$$

$$X = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \begin{bmatrix} (y - g(x, \theta)) \\ (y - g(x, \theta)) \\ \vdots \\ (y - g(x, \theta)) \end{bmatrix}$$

$$y - \cancel{y g(x, \theta)} - g(x, \theta) + \cancel{y g(x, \theta)}$$

$$\cancel{(1-g(x, \theta)) g(x, \theta)}$$

$$\Theta = \Theta + \frac{1}{N} X^T (y - g(x, \theta))$$

# Log Likelihood Derivation

$$x_i \cdot g(x, \theta) \cdot \frac{e^{-x\theta}}{1 + e^{-x\theta}}$$

$$\frac{1 + \cancel{e^{-x\theta}}}{1 + e^{-x\theta}}$$

$$- \frac{\cancel{e^{-x\theta}}}{1 + e^{-x\theta}}$$

thing

$$x_i \cdot g(x, \theta) \cdot (1 - g(x, \theta))$$

$$\frac{1}{1 + e^{-x\theta}}$$

$$1 - \text{thing} = g(x, \theta)$$

$$1 - g(x, \theta) = \text{thing}$$

# To Maximum Likelihood

$$\frac{\partial}{\partial \theta_j} \ell(y|x, \theta) = \frac{\partial}{\partial \theta_j} (y \ln(g(x, \theta)) + (1 - y) \ln(1 - g(x, \theta)))$$

- First off, a reminder...

$$\frac{\partial}{\partial x} (\ln x) = \frac{1}{x} \cdot \frac{\partial}{\partial x} (x)$$

- Therefore

$$\frac{\partial}{\partial \theta_j} \ell(y|x, \theta) = \frac{y}{g(x, \theta)} \frac{\partial}{\partial \theta_j} (g(x, \theta)) + \frac{(1 - y)}{1 - g(x, \theta)} \frac{\partial}{\partial \theta_j} (1 - g(x, \theta))$$

- But what is  $\frac{\partial}{\partial \theta_j} (g(x, \theta))$ ?

# To Maximum Likelihood

$$\frac{\partial}{\partial \theta_j} \ell(y|x, \theta) = \frac{y}{g(x, \theta)} \frac{\partial}{\partial \theta_j} (g(x, \theta)) + \frac{(1-y)}{1-g(x, \theta)} \frac{\partial}{\partial \theta_j} (1-g(x, \theta))$$

- $\frac{\partial}{\partial \theta_j} g(x, \theta) = \frac{\partial}{\partial \theta_j} \left( \frac{1}{1+e^{-x\theta}} \right) = \frac{\partial}{\partial \theta_j} (1 + e^{-x\theta})^{-1}$
- $= -1(0 - x_j e^{-x\theta}) (1 + e^{-x\theta})^{-2} = \frac{x_j e^{-x\theta}}{(1+e^{-x\theta})^2}$
- $= x_j \frac{1}{1+e^{-x\theta}} \frac{e^{-x\theta}}{1+e^{-x\theta}}$
- $= x_j g(x, \theta) (1 - g(x, \theta))$



# To Maximum Likelihood

$$\frac{\partial}{\partial \theta_j} \ell(y|x, \theta) = \frac{y}{g(x, \theta)} \frac{\partial}{\partial \theta_j} (g(x, \theta)) + \frac{(1-y)}{1-g(x, \theta)} \frac{\partial}{\partial \theta_j} (1-g(x, \theta))$$

- From the previous slide we have

$$\frac{\partial}{\partial \theta_j} g(x, \theta) = x_j g(x, \theta) (1 - g(x, \theta))$$

- Putting it all together (and simplifying) we get:

$$\frac{\partial}{\partial \theta_j} \ell(y|x, \theta) = x_j (y - g(x, \theta))$$

# To Maximum Likelihood

$$\frac{\partial}{\partial \theta_j} \ell(y|x, \theta) = x_j(y - g(x, \theta))$$

- Vectorizing this for all parameters we have

$$\frac{\partial \ell}{\partial \theta} = x^T(y - g(x, \theta))$$

- Vectorizing this to be the mean gradient over all observations we have:

$$\frac{\partial \ell}{\partial \theta} = \frac{1}{N} X^T(Y - g(X, \theta))$$

# Gradient Ascent Rule

$$\frac{\partial \ell}{\partial \theta} = \frac{1}{N} X^T (Y - g(X, \theta))$$

- We want this to go towards a maxima
- So let's update  $\theta$  as

$$\begin{aligned}\theta &:= \theta + \eta \left( \frac{\partial}{\partial \theta} \ell(Y|X, \theta) \right) \\ \theta &= \theta + \frac{\eta}{N} X^T (Y - g(X, \theta))\end{aligned}$$

- This is (almost) the same form as the least squared error for linear regression!!!!

# Non-linear decision boundaries

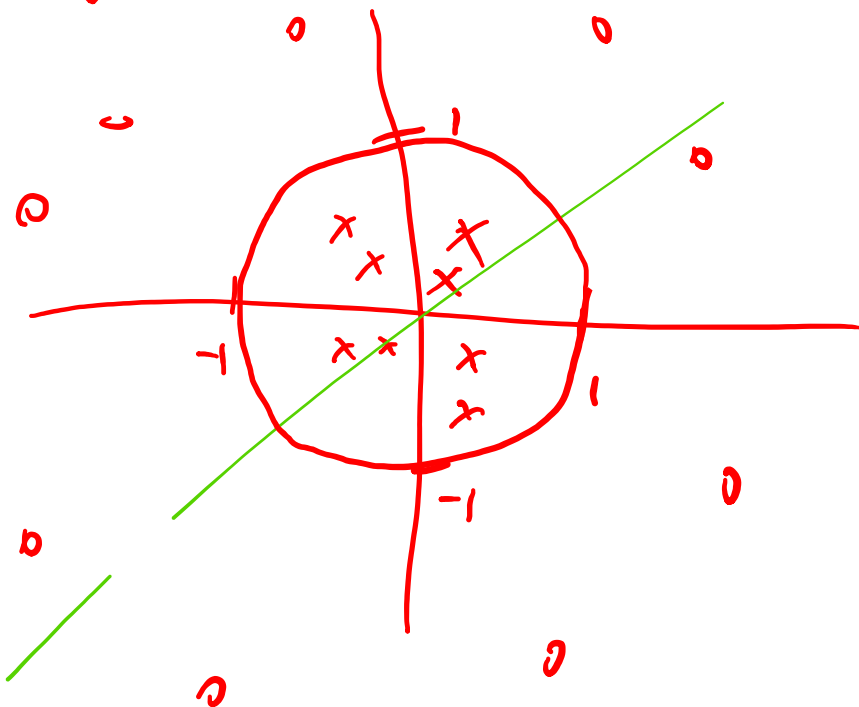


$$\theta_0 + \theta_1 x_1$$

$$\theta_0 + \theta_1 x_1 + \theta_2 x_1^2$$

...

$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2$$



$$\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} \Rightarrow -1 + x_1^2 + x_2^2 = 0$$

$$x_1^2 + x_2^2 = 1$$

$$x^2 + y^2 = 1$$