

Lecture 2: Model Building and Probability Refresher

David Carlson

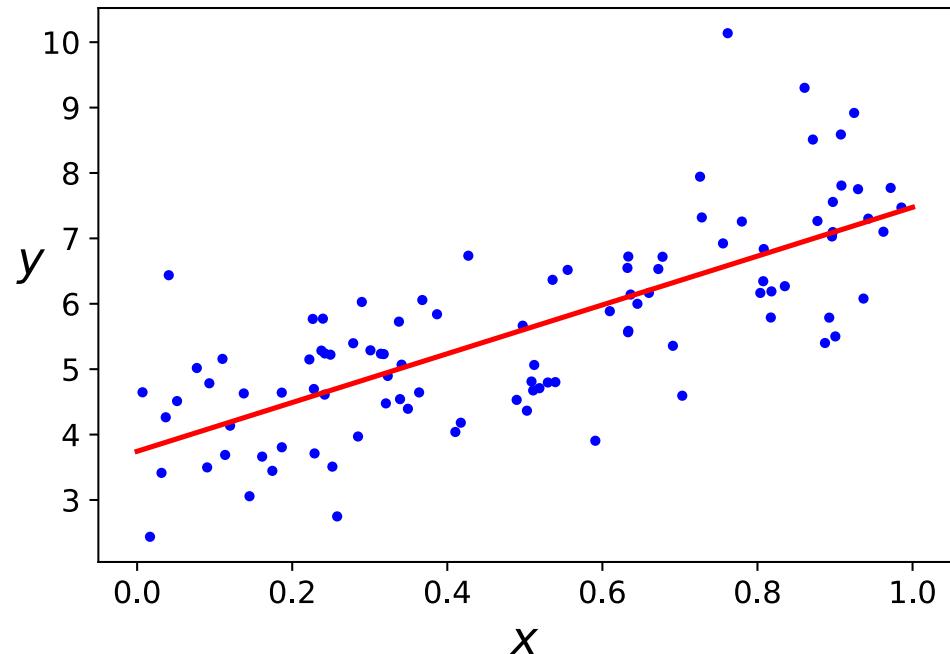
CEE 690-05: HEDS

Refresher on Predictive Modeling

- Mathematical background necessary for talking about data, statistical modeling, and machine learning
- Approach today through thinking about linear regression
- Should be review material
 - Way I talk about statistical models is often *quite* different than a statistics class

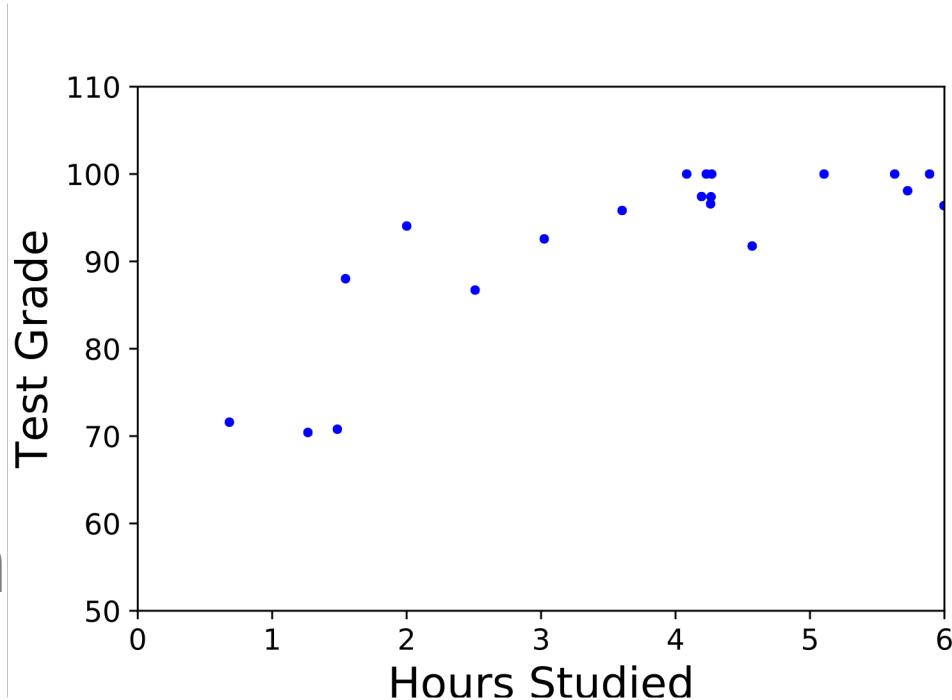
Linear Regression

- One of the earliest (and most widely used!) statistical models.
- Use to ask a variety of questions, e.g.:
 - Is smoking related to cancer?
 - Does this gene relate to height? Or to blood pressure? Or to IQ?
 - Nearly everything...
- Basis of many idea in statistics and the *underlying* (but not complete) model in many research areas, such as genomics
- Everyone should have seen linear regression before



Predicting Test Scores

- Suppose that we want to predict test scores from time spent studying
- How would we set up linear regression to do that?



What are we trying to do?



x , data for a subject (i.e.
amount of time studying)

y , associated test score

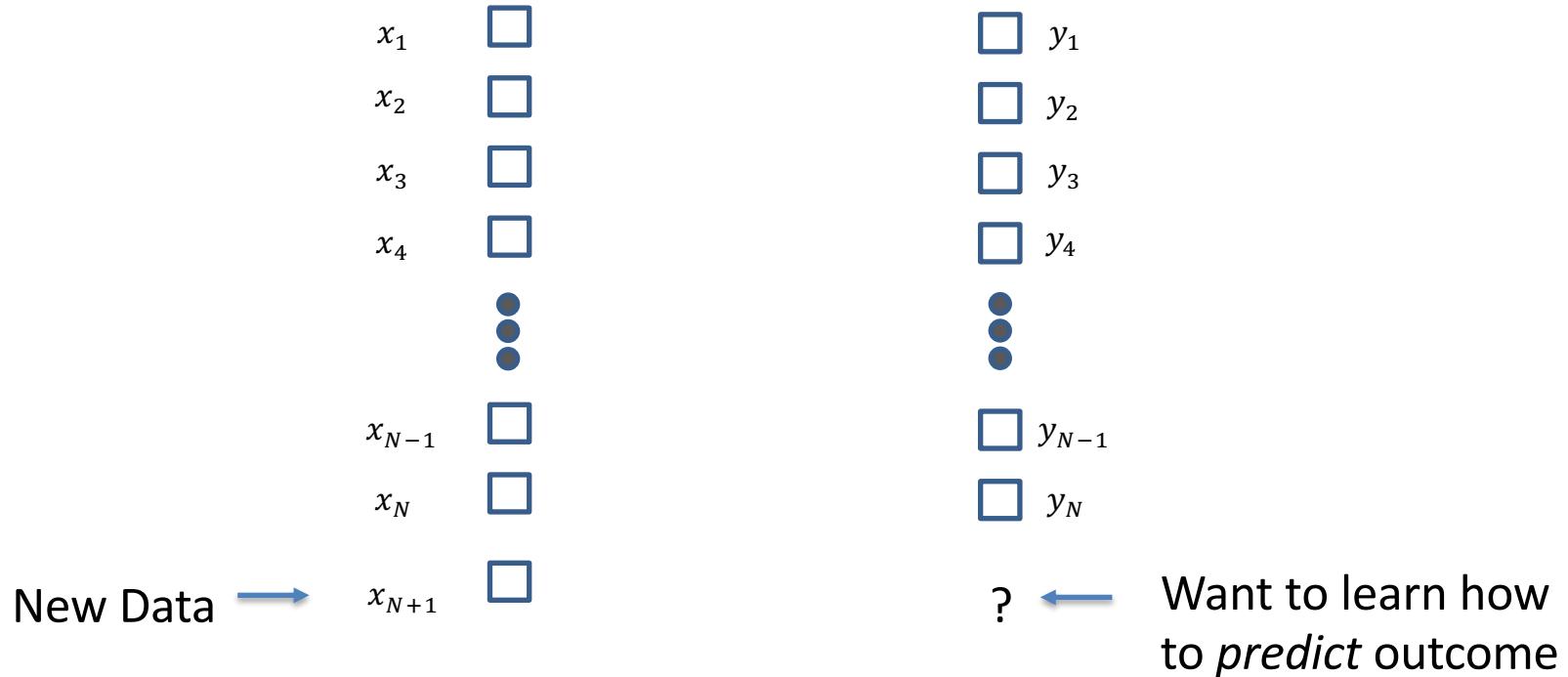
End goal: *predict* y from x

Training Dataset/ Historical Examples

x_1	<input type="checkbox"/>	<input type="checkbox"/>	y_1
x_2	<input type="checkbox"/>	<input type="checkbox"/>	y_2
x_3	<input type="checkbox"/>	<input type="checkbox"/>	y_3
x_4	<input type="checkbox"/>	<input type="checkbox"/>	y_4
	⋮	⋮	
x_{N-1}	<input type="checkbox"/>	<input type="checkbox"/>	y_{N-1}
x_N	<input type="checkbox"/>	<input type="checkbox"/>	y_N

Time spend studying Test Score

Training Dataset/ Historical Examples



Linear Regression

- Linear regression chooses the function:
$$\hat{y} = f(x; \beta_0, \beta_1) = \beta_0 + \beta_1 x,$$
- Where β_0 and β_1 are constant parameters and x is a variable.
- Need to choose some way of “learning” what β_0 and β_1 are for this data problem (i.e. statistical inference)

Minimizing (Mean) Squared Error

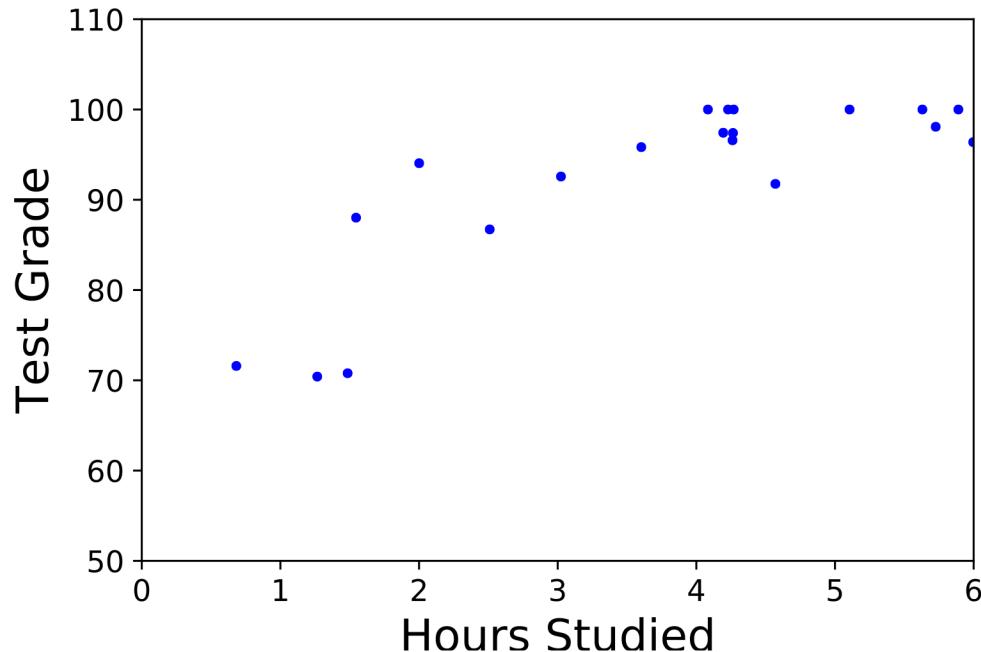
- In linear regression, we nearly* always choose to minimize the squared error over all observed data points:

$$\beta_0^*, \beta_1^* = \arg \min_{\beta_0, \beta_1, \hat{y}=\beta_0 + \beta_1 x} \sum_{i=1}^N (y_i - \hat{y})^2$$

- Why is this the square? (will come back to)
- How do we actually solve this? (will come back to)

Fitting a “simple” linear regression

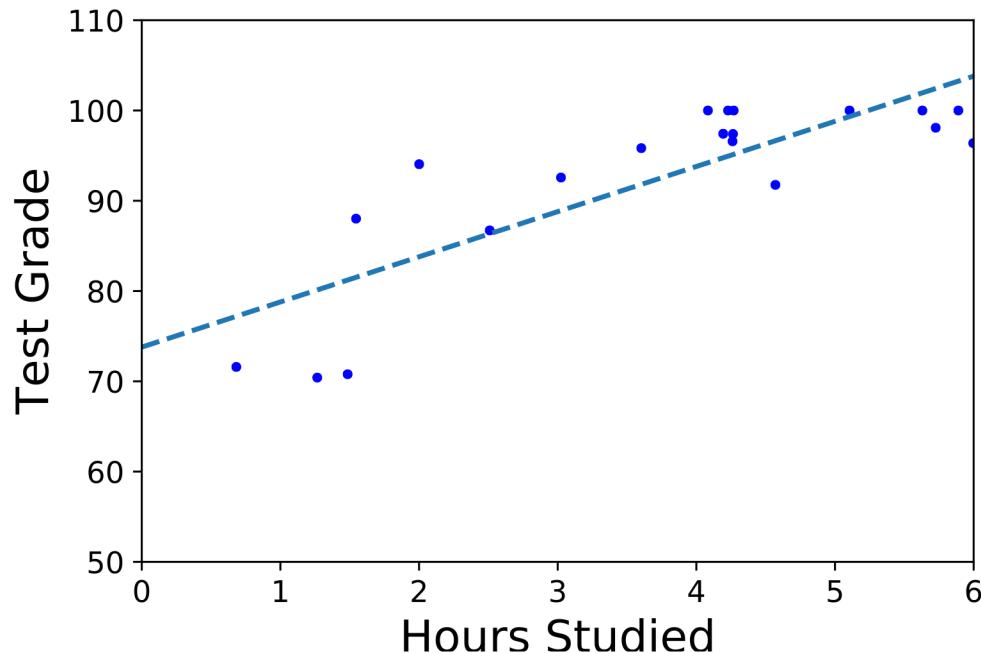
Suppose we have the data to the right, what does our linear regression look like here?



Fitting a “simple” linear regression

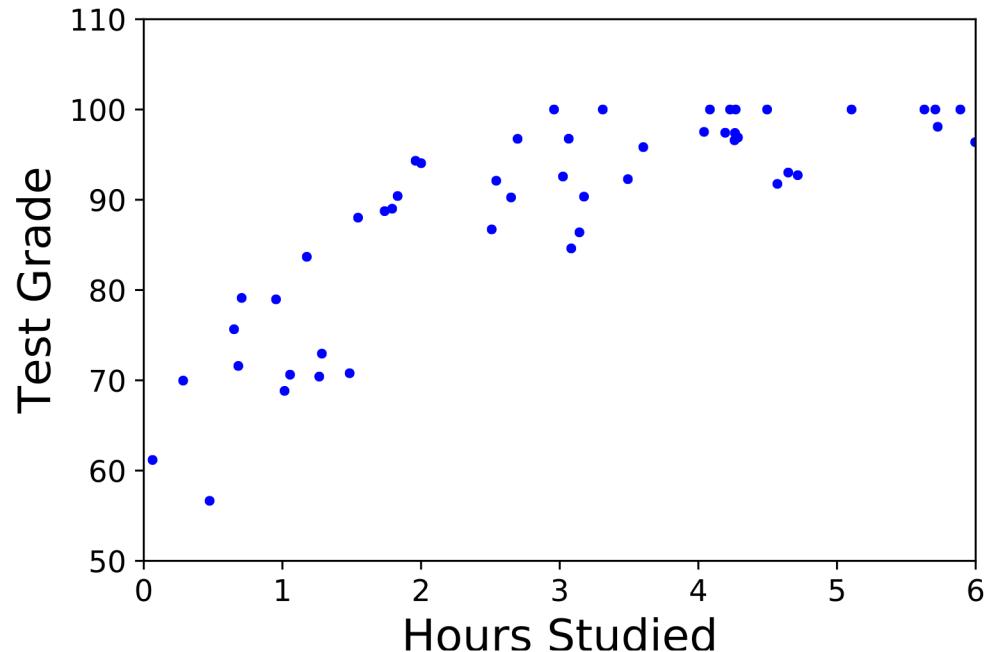
Suppose we have the data to the right, what does our linear regression look like here?

The linear regression fit does a “reasonable” job predicting scores here.



Can keep updating our model as we see more examples

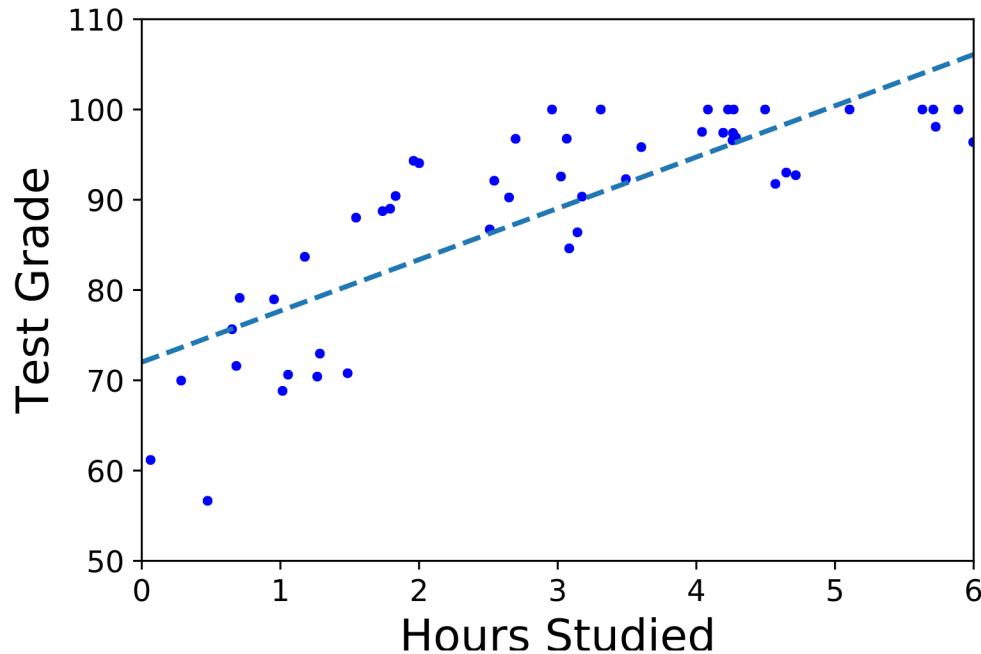
As data keeps coming in, we can use our additional examples (i.e. past experiences) to make improved predictions about the future.



Can keep updating our model as we see more examples

As data keeps coming in, we can use our additional examples (i.e. past experiences) to make improved predictions about the future.

Fitting linear regression again...

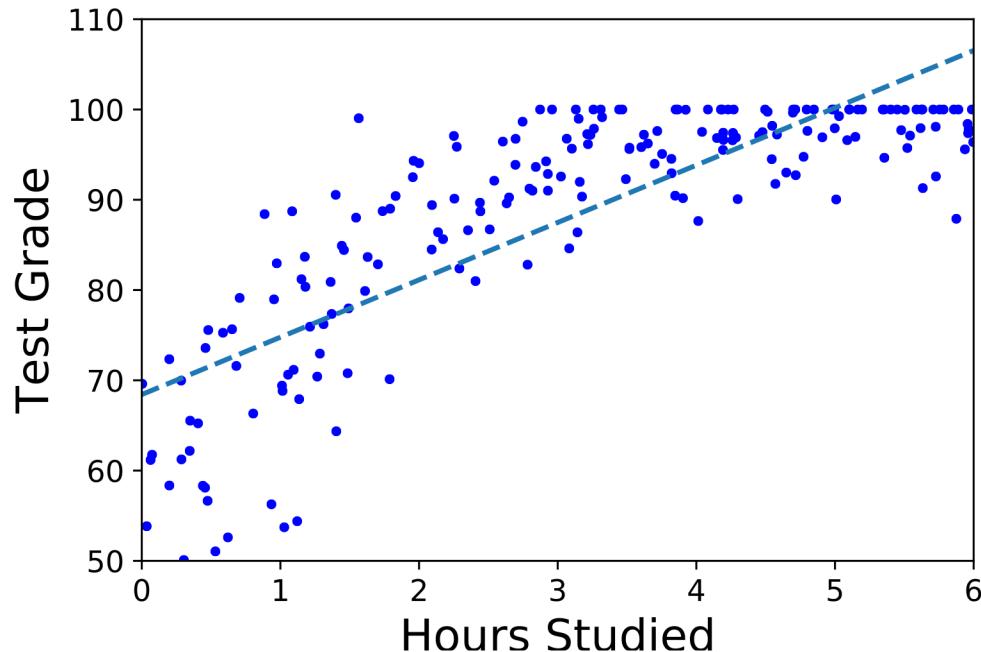


Can keep updating our model as we see more examples

As data keeps coming in, we can use our additional examples (i.e. past experiences) to make improved predictions about the future.

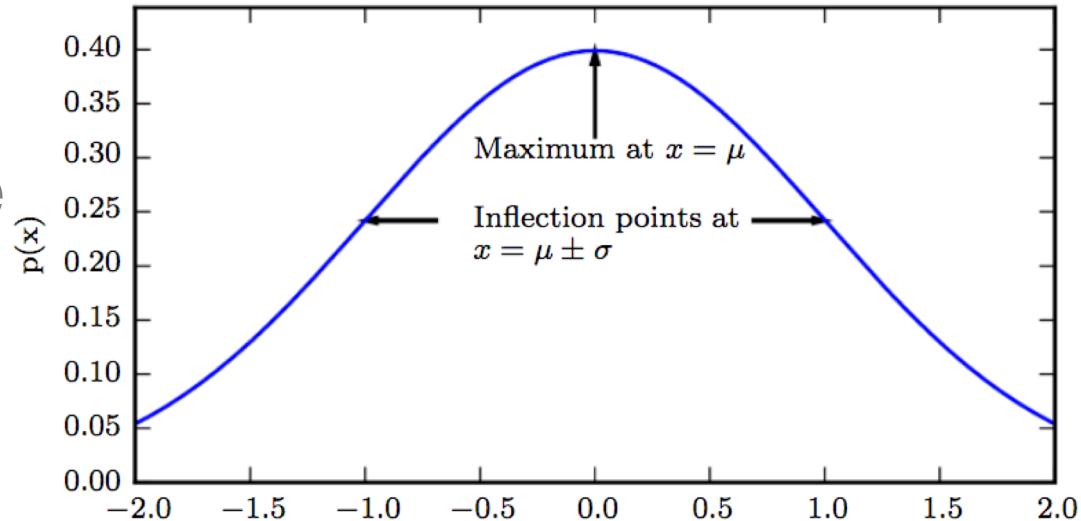
Adding much more data...

Is this still a reasonable data fit?



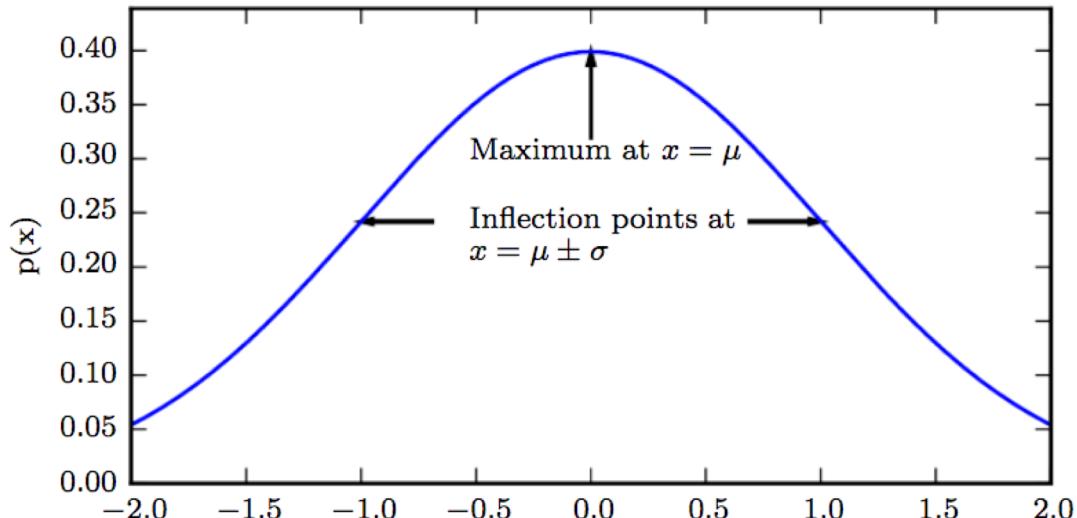
Why do we used squared error?

- It all has to do with our normal/Gaussian/ bell curve distribution.
- Often we say that errors are “normally distributed”



Quick Refresh on Normal Distribution

- A variable with a normal distribution is defined:
 $x \sim N(\mu, \sigma^2)$
- With probability density function:
 $p(x; \mu, \sigma^2)$
 $= \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)$
- The “standard normal” is shown at the right with $\mu = 0$ and $\sigma^2 = 1$



Linear Regression

- We often implicitly consider the statistical model

$$y = \beta_0 + \beta_1 x + \epsilon$$

- The “error” ϵ is normally distributed, where

$$\epsilon = y - \beta_0 + \beta_1 x$$

Maximum Probability

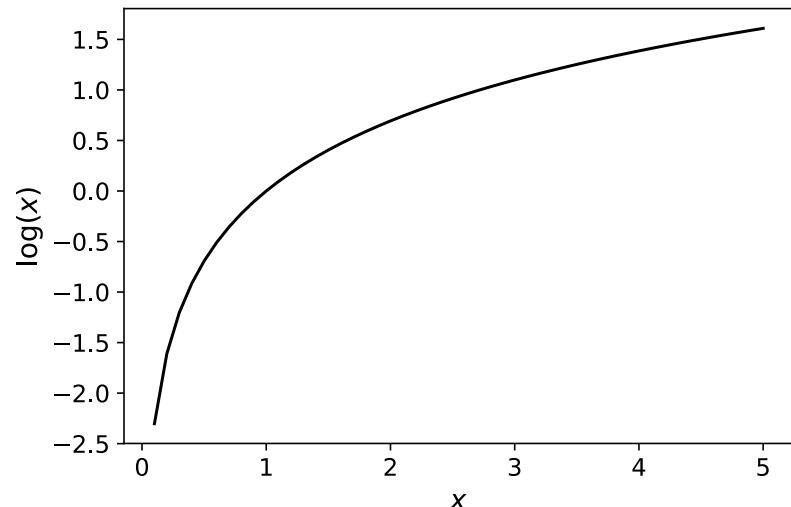
- We want to *maximize* statistical probability of our data
 - (i.e. our observations were as *likely* as possible under our model)
- This means

$$\beta_0^*, \beta_1^* = \arg \max_{\beta_0, \beta_1, \epsilon = y - \beta_0 - \beta_1 x} \prod_{i=1}^N p(\epsilon_i)$$

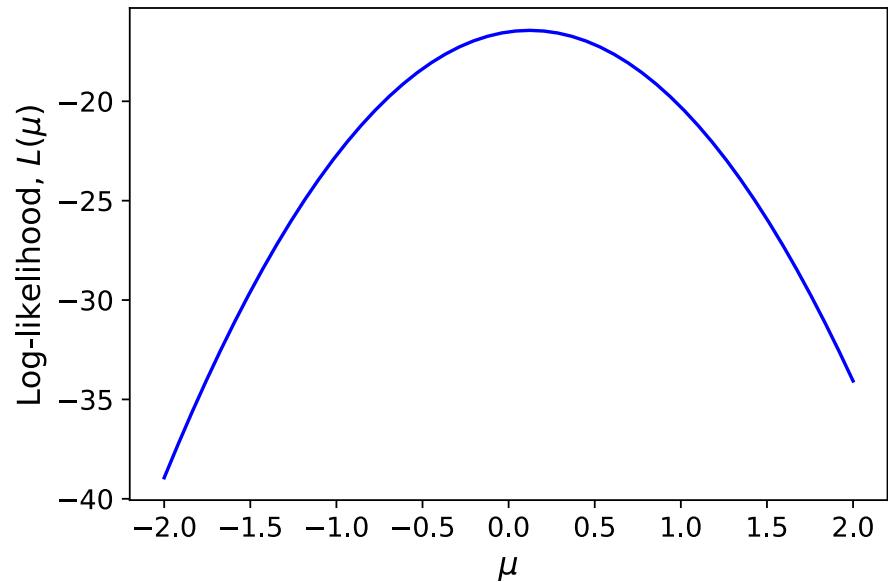
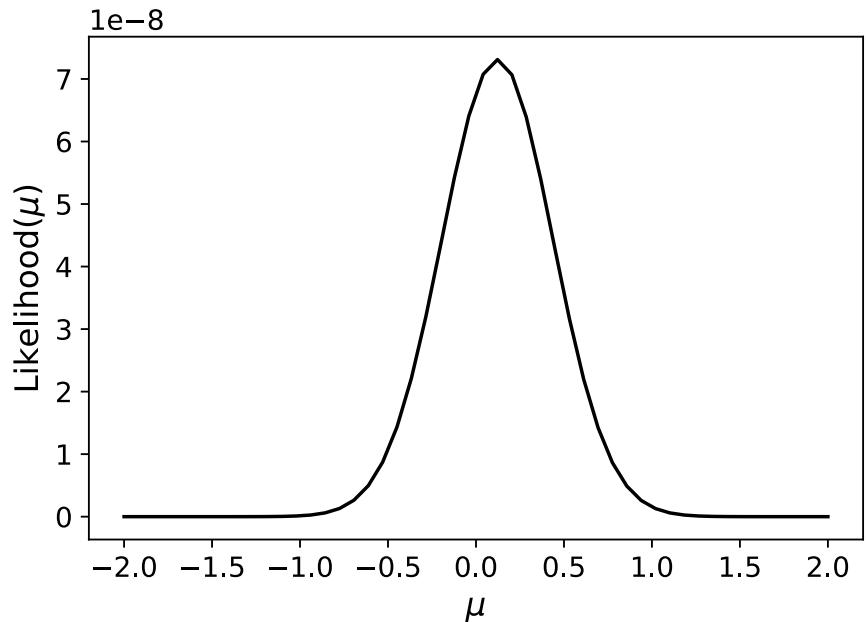
$$\beta_0^*, \beta_1^* = \arg \max_{\beta_0, \beta_1, \epsilon = y - \beta_0 - \beta_1 x} \prod_{i=1}^N N(\epsilon_i; 0, \sigma^2)$$

Monotonic transformations

- Dealing with products that arise in log-likelihoods is frustrating mathematically, instead we use the *log* function to transform the likelihood.
- Because the function is monotonic, the location of the maximum point does not change, which we'll visualize on the next slide.



Normal (Log-)Likelihood



Maximum Log-Likelihood

- This means

$$\beta_0^*, \beta_1^* = \arg \max_{\beta_0, \beta_1, \epsilon = y - \beta_0 - \beta_1 x} \log \prod_{i=1}^N p(\epsilon_i)$$

$$\beta_0^*, \beta_1^* = \arg \max_{\beta_0, \beta_1, \epsilon = y - \beta_0 - \beta_1 x} \sum_{i=1}^N \log p(\epsilon_i)$$

$$\beta_0^*, \beta_1^* = \arg \min_{\beta_0, \beta_1, \hat{y} = \beta_0 + \beta_1 x} const_1 + const_2 \sum_{i=1}^N (y_i - \hat{y})^2$$

Some considerations

- We assume that we can solve the previous math equation
- Often, we have many features of a data (e.g hours studied, grade on the last test, grade on the last homework,...)
- Often, we care about predicting other types of outcomes, such as binary (positive/negative)

Learning a Predictive Model with Many Features



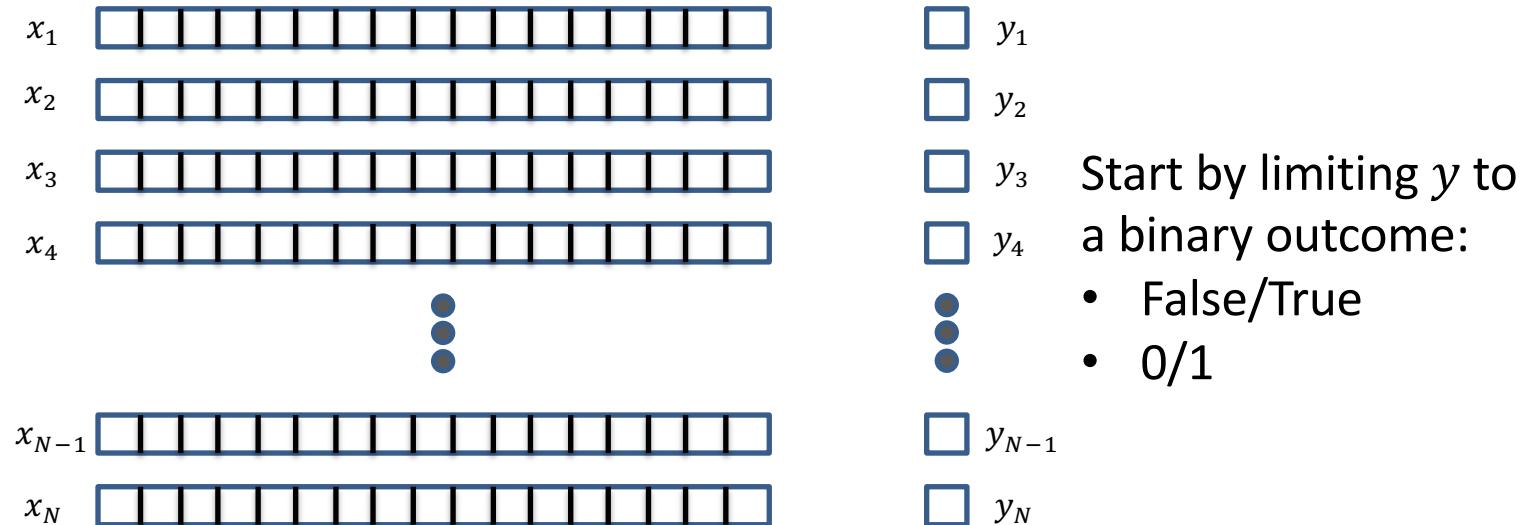
x , data/features for
a subject



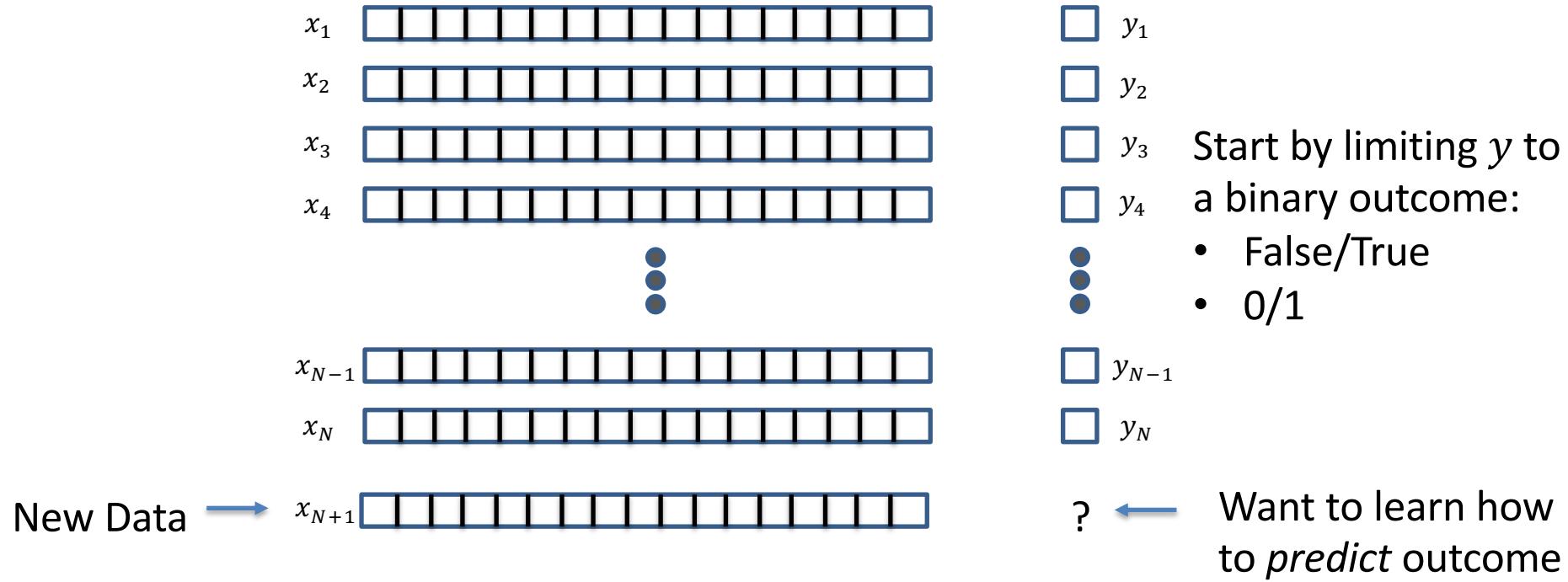
y , associated label 0/1

End goal: *predict* y from x

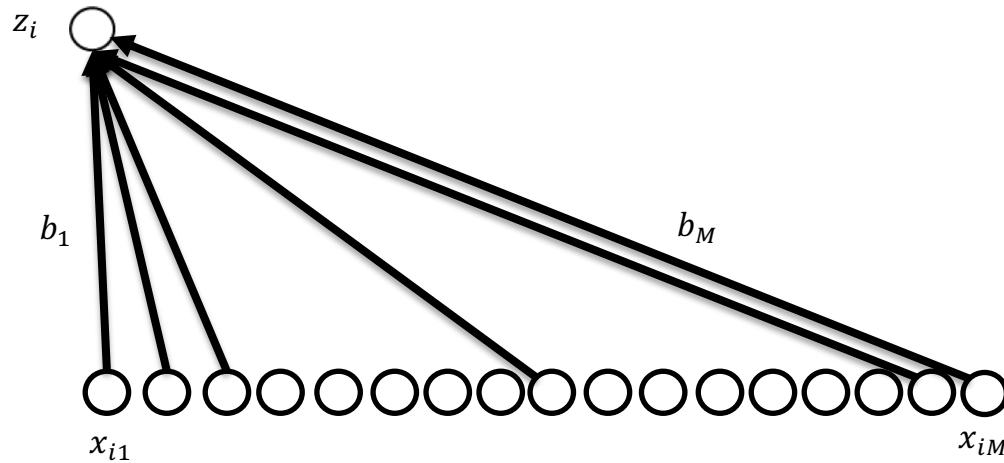
Training Set (Historical Data)



Making Predictions



Linear Predictive Model



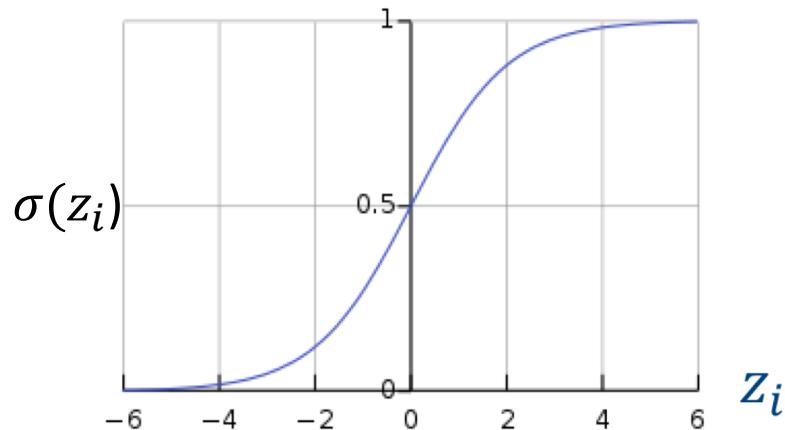
$$z_i = (b_1 \times x_{i1}) + (b_2 \times x_{i2}) + \cdots + (b_M \times x_{iM})$$

Convert to a Probability

$$z_i = (b_1 \times x_{i1}) + (b_2 \times x_{i2}) + \cdots + (b_M \times x_{iM}) + b_0$$

$$p(y_i = 1|x_i) = \sigma(z_i)$$

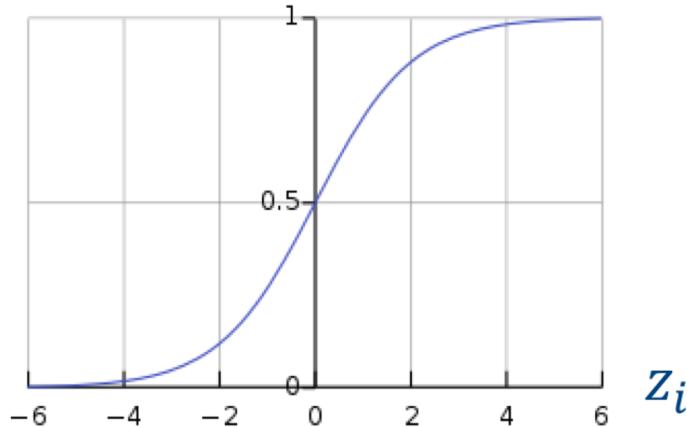
Extra Constant



Convert to a Probability

$$z_i = (b_1 \times x_{i1}) + (b_2 \times x_{i2}) + \cdots + (b_M \times x_{iM}) + b_0$$

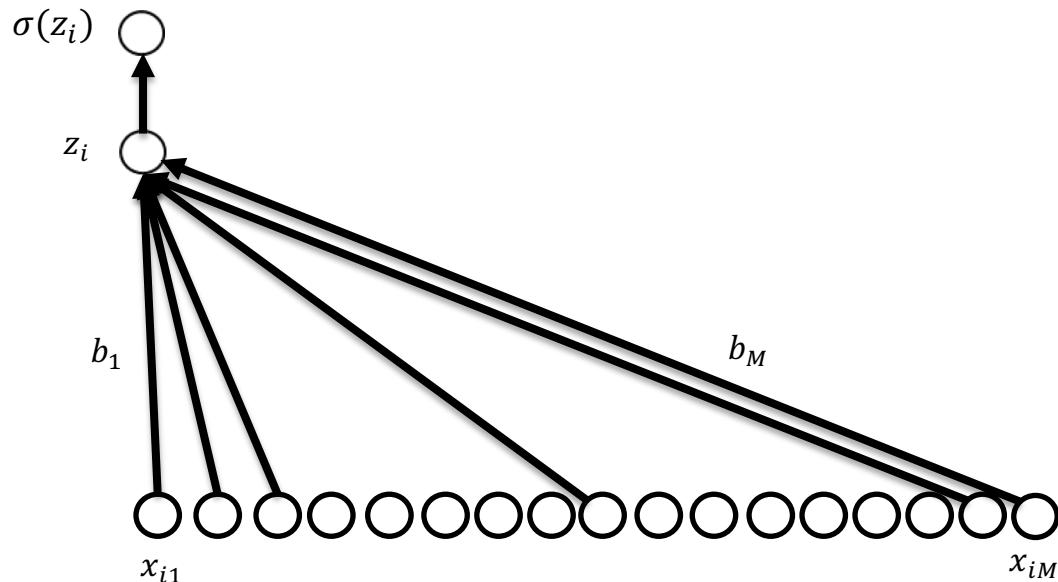
$$p(y_i = 1|x_i) = \sigma(z_i) = \frac{\exp(z_i)}{1+\exp(z_i)}$$



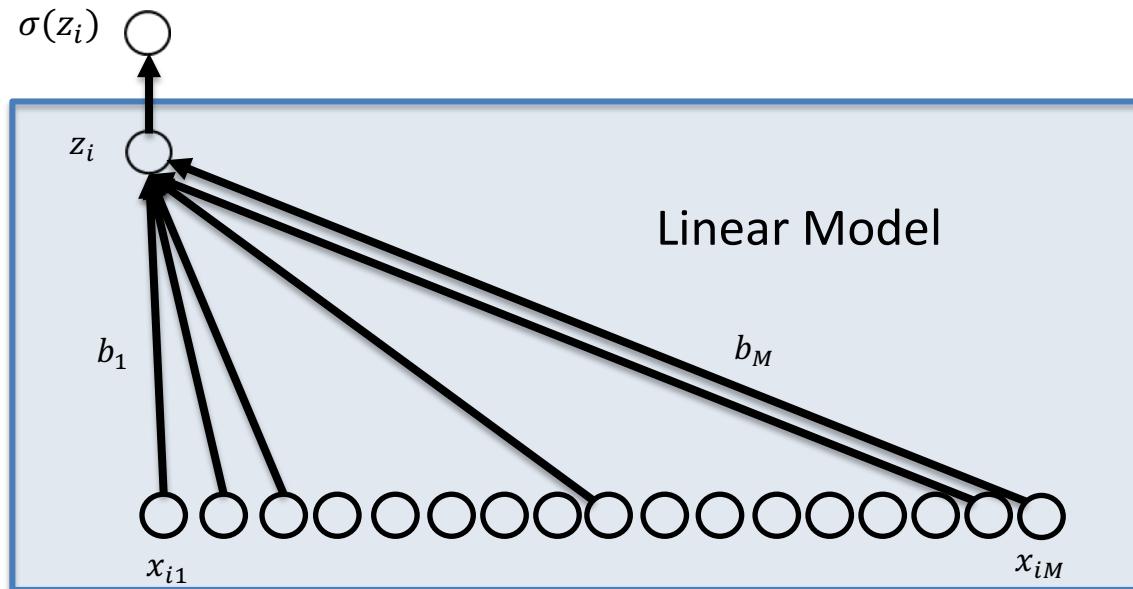
- Large and positive z_i indicates that event $y_i = 1$ is likely

- Large and negative z_i indicates that event $y_i = 0$ is likely

Logistic Regression

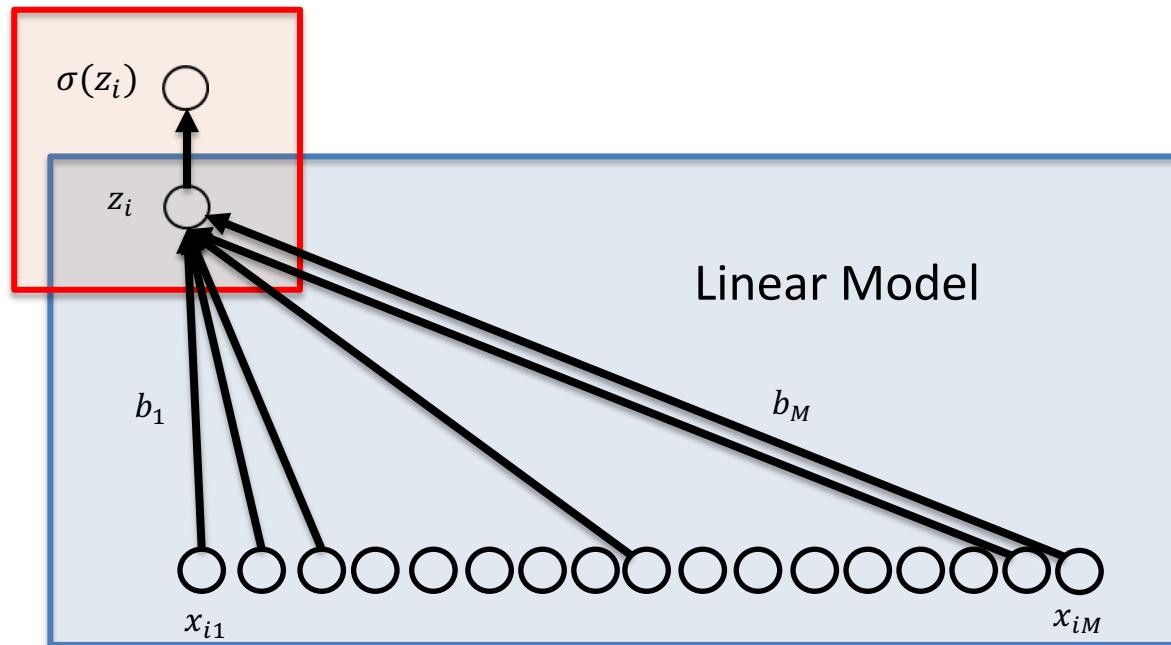


Logistic Regression



Logistic Regression

Convert to
Probability



What do the parameters and model mean?

AN EXAMPLE

Example

- Outcome:
 - $y_i = 1$, it rains on day i ;
 - $y_i = 0$, it does not rain on day i
- Features:
 - On day i what is the *{cloud cover, humidity, temperature, air pressure, ...}*



x_i , features for day i

y_i , did it rain on day i

Example

- **Outcome:** $y_i = 1$, it rains on day i ; $y_i = 0$, it does not rain on day i
- **Features:** On day i what is the {1: *cloud cover*, 2: *humidity*, 3: *temperature*, ... }

$$z_i = (b_1 \times x_{i1}) + (b_2 \times x_{i2}) + \cdots + (b_M \times x_{iM}) + b_0$$

Cloud Cover Humidity

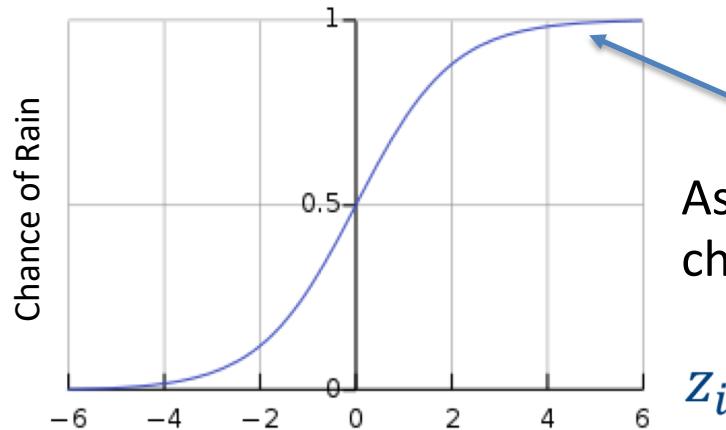
- If cloud cover is positively related to rainfall, b_1 should be positive

Impact on the Sigmoid Function

$$z_i = (b_1 \times x_{i1}) + (b_2 \times x_{i2}) + \dots + (b_M \times x_{iM}) + b_0$$

Cloud Cover *Humidity*

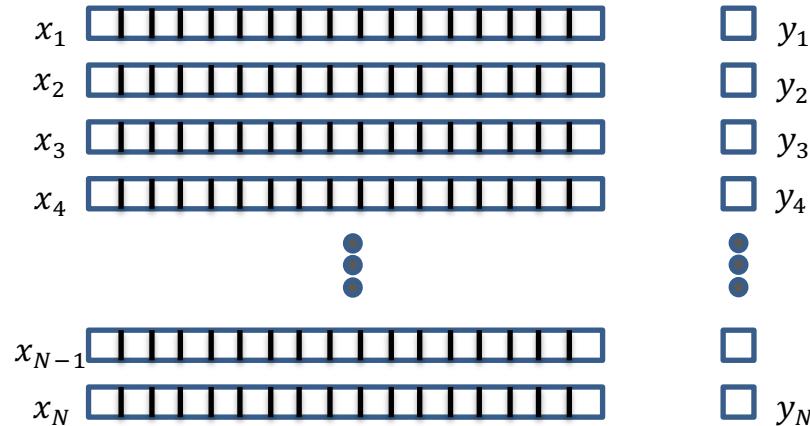
$$p(y_i = 1|x_i) = \sigma(z_i)$$



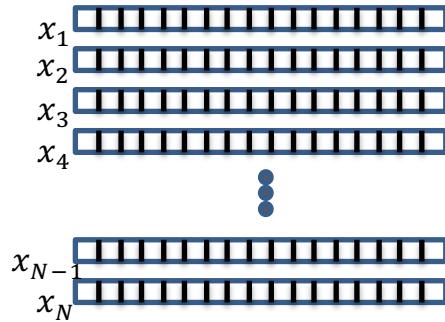
As the value z_i increases, the chance of rain increases

Building the Training Set

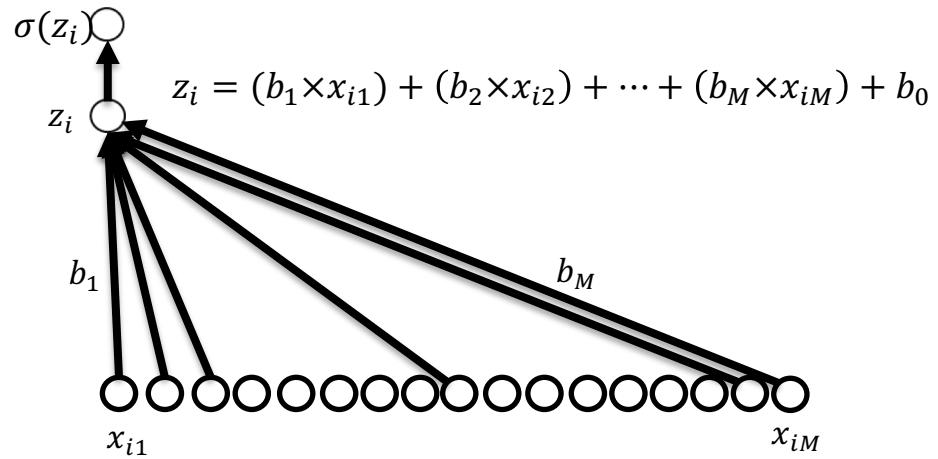
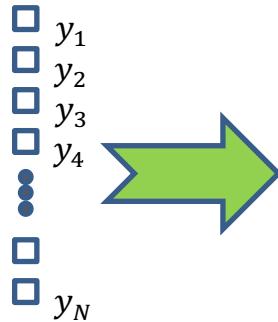
- Need to learn the parameters
- Requires *training data*
- Record data from N days
 - Capture features: {*cloud cover, humidity, temperature, ...*}
 - Did it rain?



Learning Model Parameters



Training Set

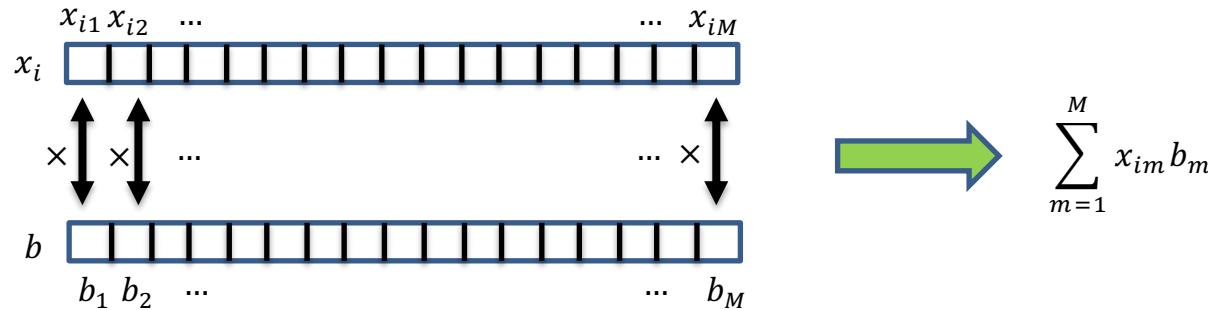


Logistic Regression Model (or “Network”)

A large green downward-pointing arrow labeled "Learned Parameters" with the expression (b_0, b_1, \dots, b_N) below it, indicating the result of the learning process.

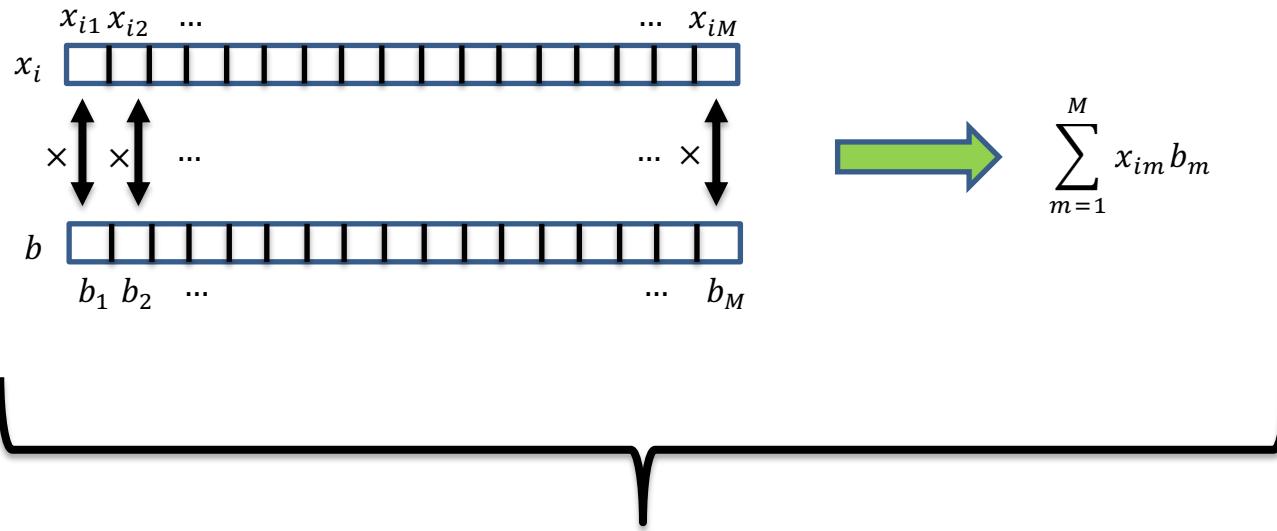
Interpretation of Logistic Regression

$$z_i = (b_1 \times x_{i1}) + (b_2 \times x_{i2}) + \cdots + (b_M \times x_{iM})$$



Interpretation of Logistic Regression

$$z_i = (b_1 \times x_{i1}) + (b_2 \times x_{i2}) + \cdots + (b_M \times x_{iM})$$



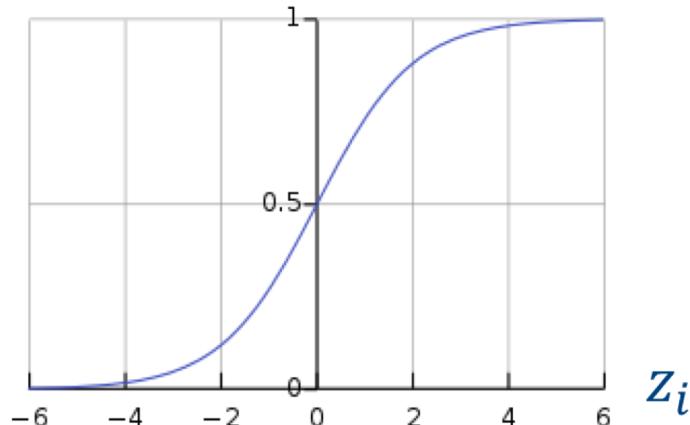
Compact Notation: $x_i \odot b$ (or “inner product”)

Interpretation of Logistic Regression

$$z_i = b_0 + (b_1 \times x_{i1}) + (b_2 \times x_{i2}) + \cdots + (b_M \times x_{iM})$$

$$= b_0 + x_i \odot b$$

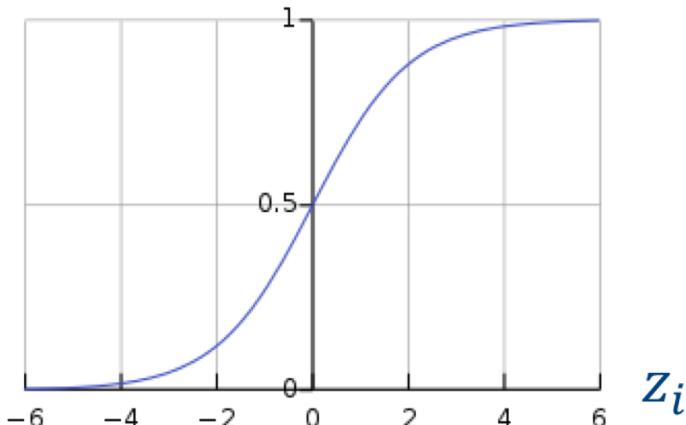
$$p(y_i = 1|x_i) = \sigma(z_i)$$



Interpretation of Logistic Regression

$$\begin{aligned}z_i &= b_0 + (b_1 \times x_{i1}) + (b_2 \times x_{i2}) + \cdots + (b_M \times x_{iM}) \\&= b_0 + x_i \odot b\end{aligned}$$

$$p(y_i = 1|x_i) = \sigma(z_i)$$



- May think of vector b as a template or filter (will visualize to make clear)
- If x_i is aligned/matched with b , then $x_i \odot b$ will be large
- The parameter b_0 is a bias to correct for class prevalences

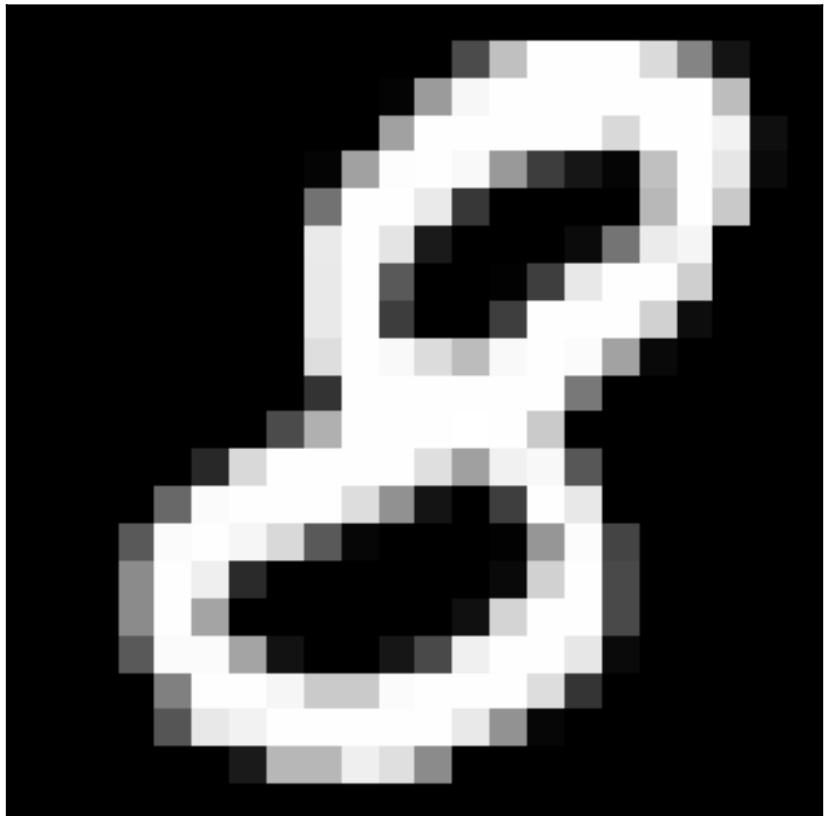
RECOGNIZING NUMERICAL DIGITS

The MNIST Dataset

- The Modified National Institute of Standards and Technology (MNIST) contains pictures of handwritten digits (0,1,2,...)
- Want to be able to tell what digit each image is (e.g. optical character recognition)

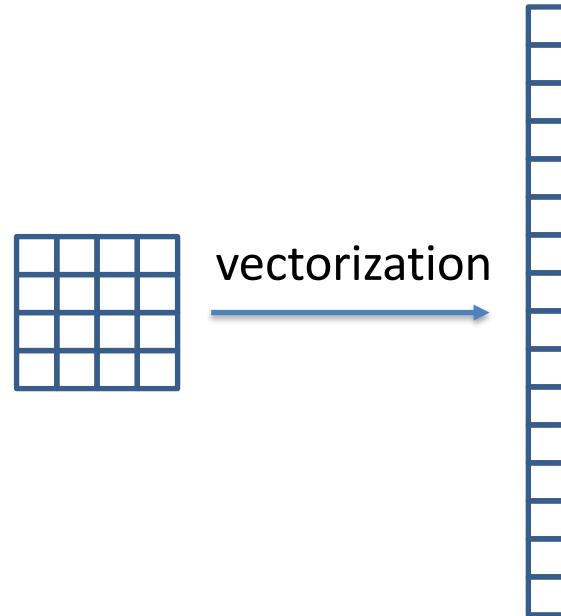


Images are Encoded as Numbers



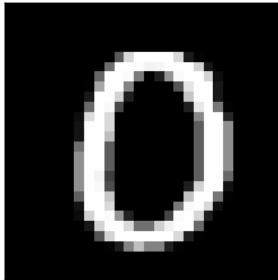
Vectorization

- We will start talking about deep learning *without* using the structure of the image
- We will return to this in a future lecture
- To convert an image into an unstructured set of numbers, we *vectorize it*

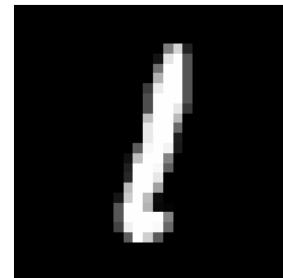
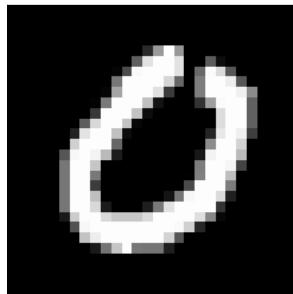
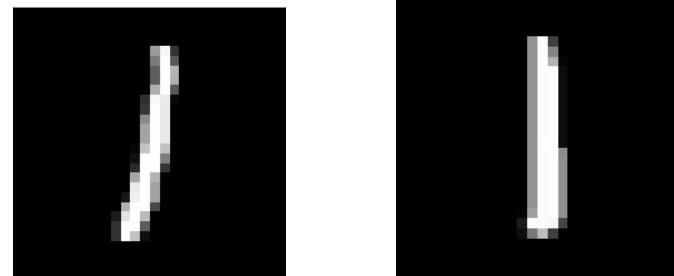


Start With The Binary Case

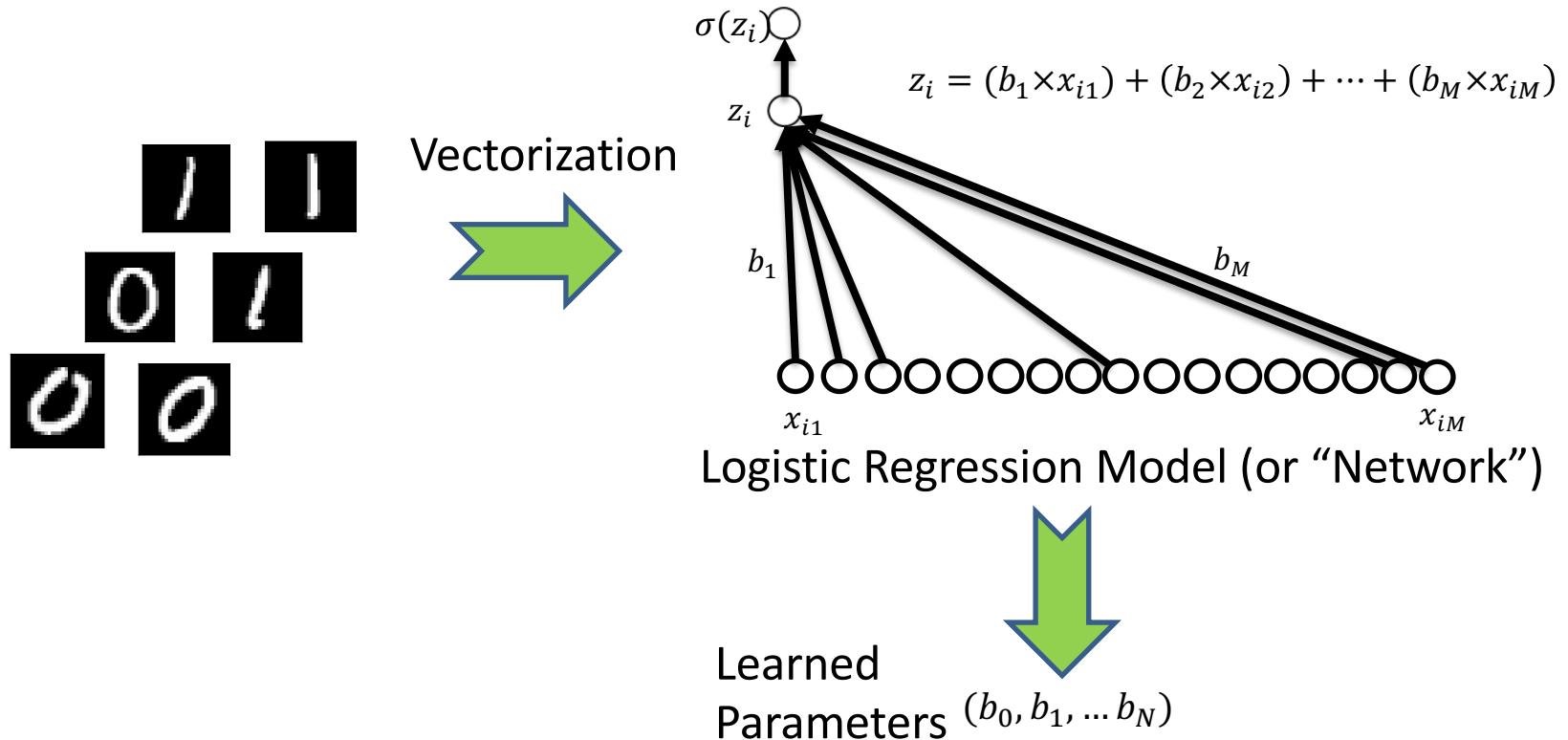
Zeros



Ones

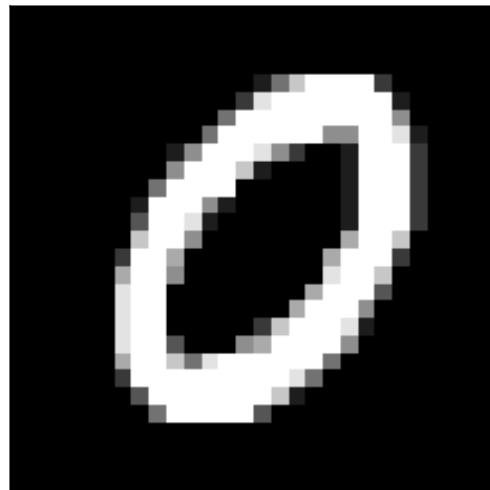


Learning on MNIST

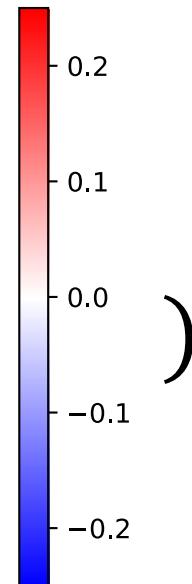
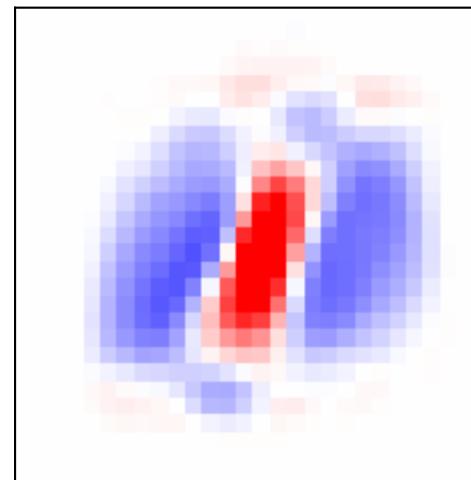


Zooming in on 0/1

$\sigma($



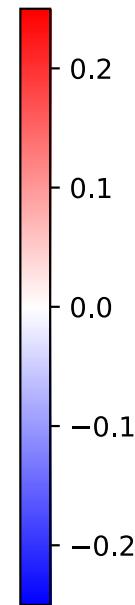
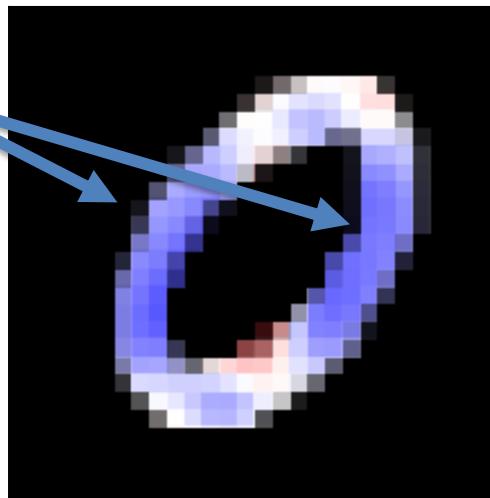
\times



Zooming in on 0/1

Negative Sections

$\sigma($

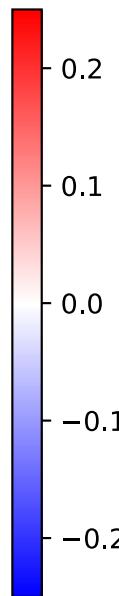
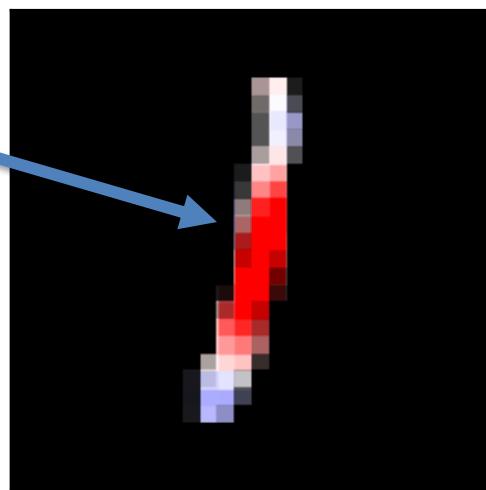


) = 0.006

Zooming in on 0/1

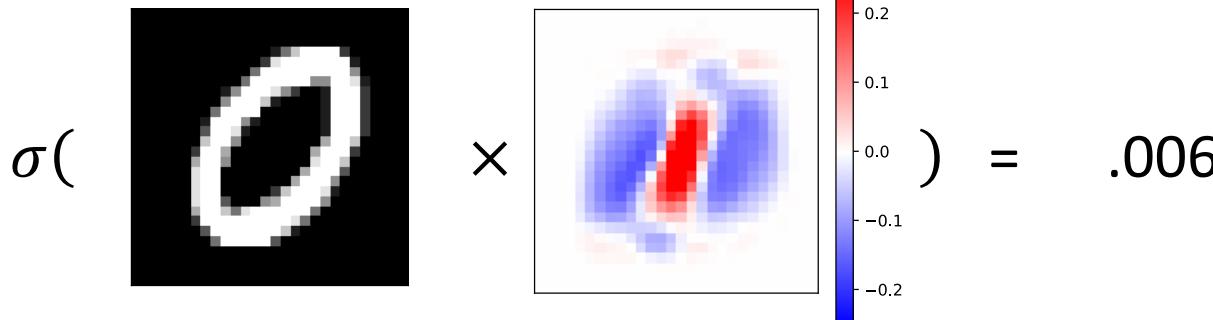
Positive Section

$\sigma($

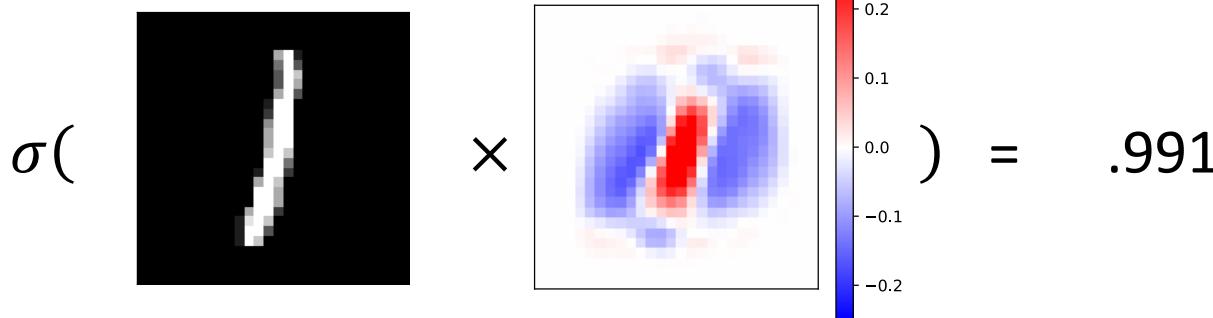


) = .991

Learned Weights for 0/1



We think that
this is a “zero”
(.6% chance it
is a “one”)



We think that
this is a “one”
(99.1% chance
it is a “one”)

How do we determine the parameters here?

- Same as before, we're going to use the log-likelihood
- This time, we're going to use the Bernoulli (e.g. binary, binomial)
$$p(y = 1; r) = r, p(y = 0; r) = (1 - r)$$
- This is the same as a weighted coin flip

How do we determine the parameters here?

- In logistic regression, we have

$$p(y = 1; r) = r,$$

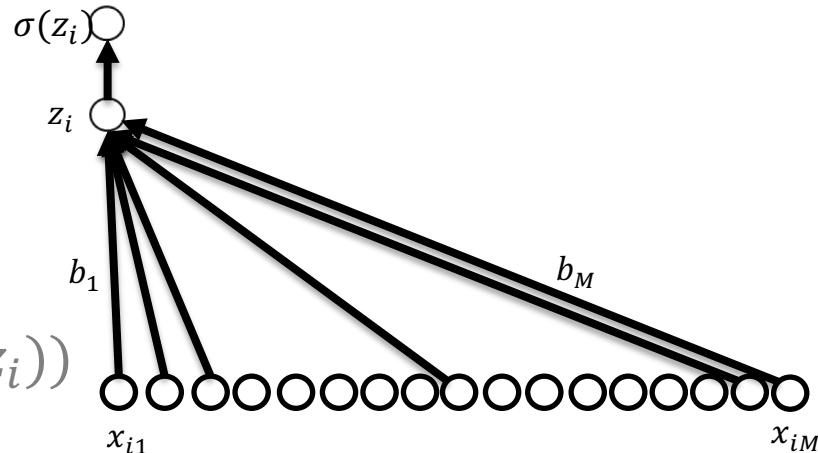
$$p(y = 0, r) = (1 - r)$$

$$r = \sigma(b_0 + x_i \odot b)$$

- Maximize:

$$b_0^*, b^*$$

$$= \arg \max_{b_0, b, z_i = b_0 + x_i \odot b} \sum_{i=1}^N \log p(y_i; \sigma(z_i))$$

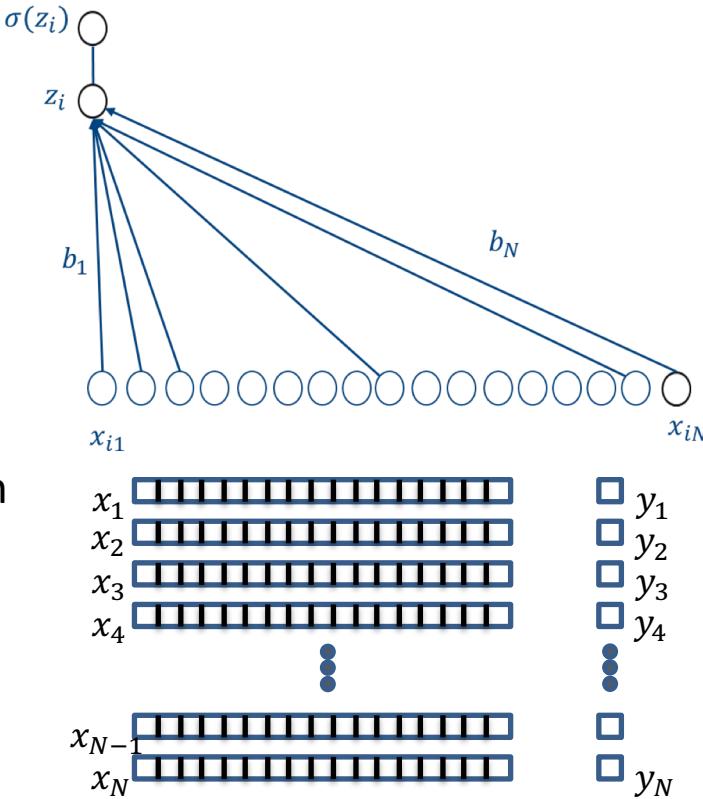


“Empirical Risk Minimization”

- A *loss function* $\ell(\text{true}, \text{prediction})$ defines a penalty for poor predictions
- Want to minimize average loss
- Mathematically, this can be stated as

$$\mathbf{b}^* = \arg \min_{\mathbf{b}} \frac{1}{N} \sum_i^N \ell(y_i, \sigma(z_i))$$

↑ ↑
True Label Loss function
 Guess

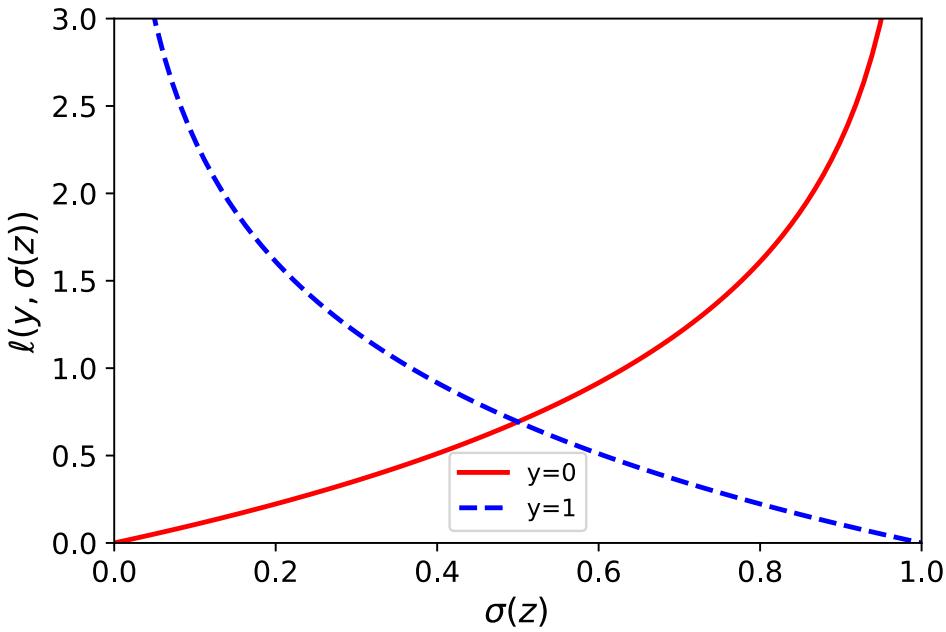


What is the Loss Function?

- Define $\sigma(z_i)$ as the predicted probability
- y_i is our true label
- Can be viewed as the negative log-likelihood
$$\ell(y_i, \sigma(z_i)) = -\log p(y_i | \sigma(z_i))$$
- Specific mathematical form is:
$$\ell(y, \sigma(z)) = -y \log \sigma(z) - (1 - y) \log(1 - \sigma(z))$$

Visualization of Logistic Loss Function

- The logistic/cross-entropy loss is:
$$\ell(y, \sigma(z)) = -y \log \sigma(z) - (1 - y) \log(1 - \sigma(z))$$
- Dashed blue line shows loss when the true prediction is positive
- When we give 100% of a 1, we pay no penalty
- If we are less confident, we pay a small penalty
- If we are overconfident and wrong, we pay a large penalty

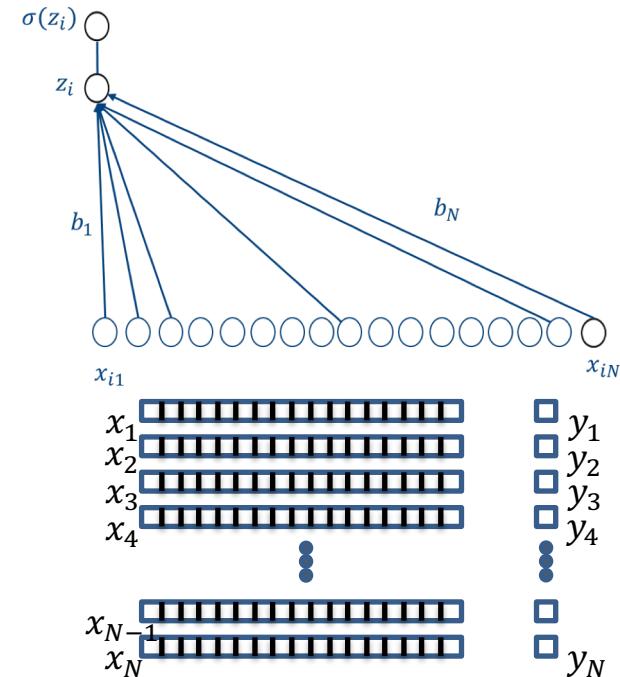


Optimization Goal for Binary Classification

- The optimization goal is to minimize the average loss

$$\mathbf{b}^* = \arg \min_{\mathbf{b}} \frac{1}{N} \sum_i^N \ell(y_i, \sigma(z_i))$$

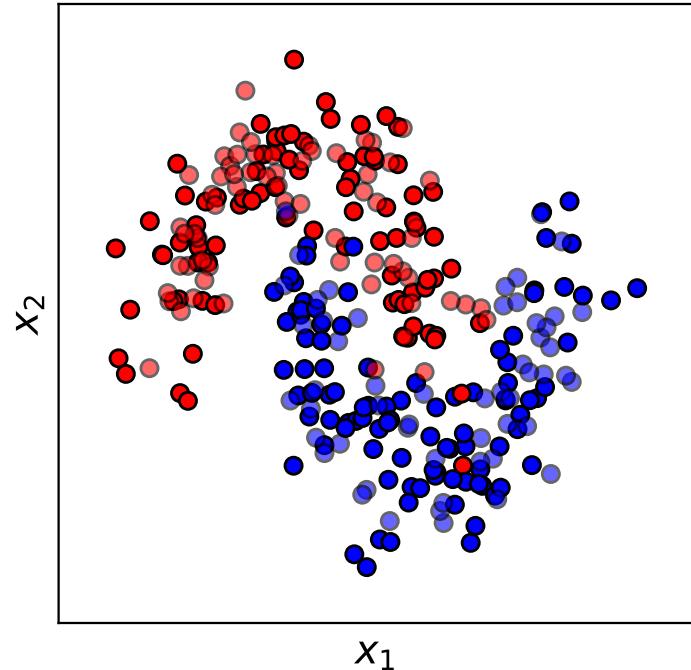
- For binary (0/1) problems, the logistic or cross-entropy loss is
$$\ell(y, \sigma(z)) = -y \log \sigma(z) - (1 - y) \log(1 - \sigma(z))$$



LIMITATIONS, LOOKING FORWARD

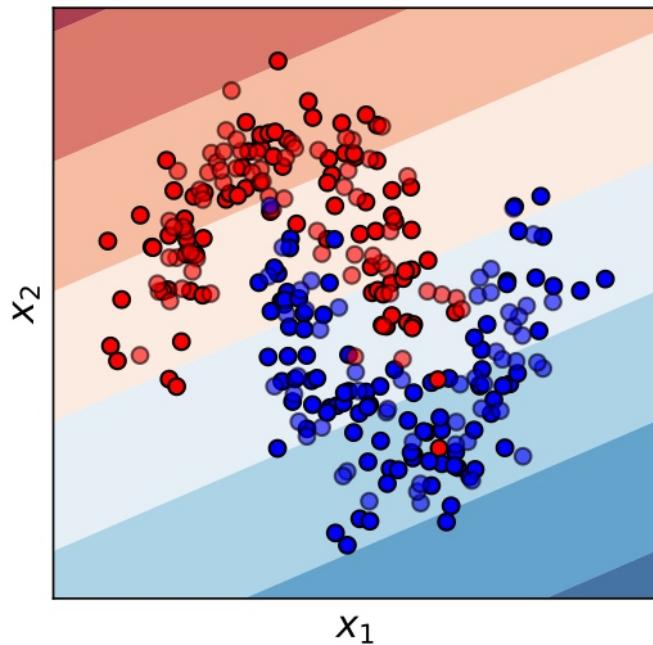
Logistic Regression is a “Linear” Classifier

- Also referred to as a generalized linear model
- Can only split data by linear trends



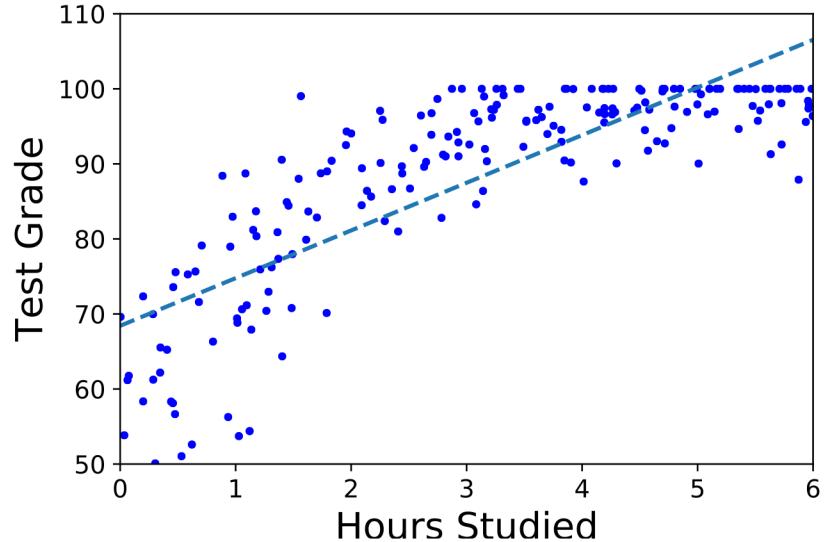
Limitations of Linear Classifiers

- Linear classifiers can only represent limited relationships
- Often want to use a classifier that can handle non-linearities.
- Many different ways of creating non-linear classifiers!



Can we create a non-linear model for test scores?

- Want to capture the effect of “diminishing returns” in the data
- Can we construct automatic ways of “learning” this?



Looking Forward

- These are two initial techniques for making predictions from data
- Will start talking about machine learning algorithms for making sense of the data