

INTRO to DATA SCIENCE

LECTURE 9: PROBABILITY & LOGISTIC REGR.

LAST TIME

LINEAR REGRESSION

POLYNOMIAL REGRESSION

REGULARIZATION

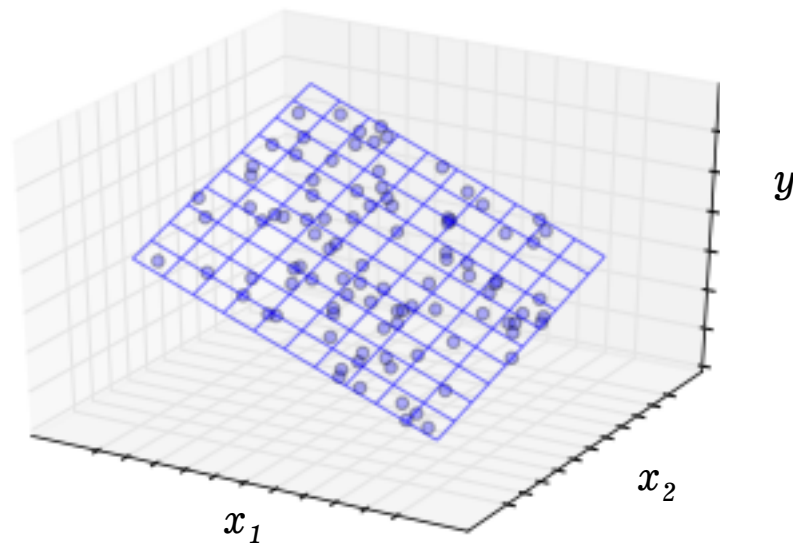
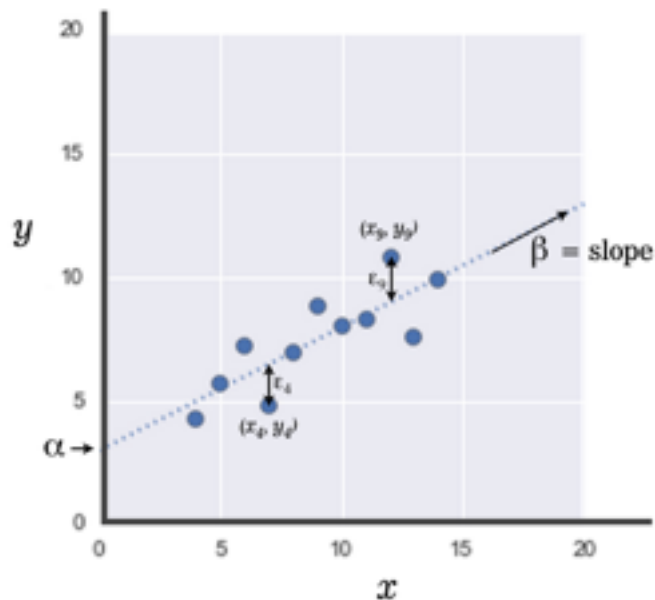
FEATURE CREATION

...

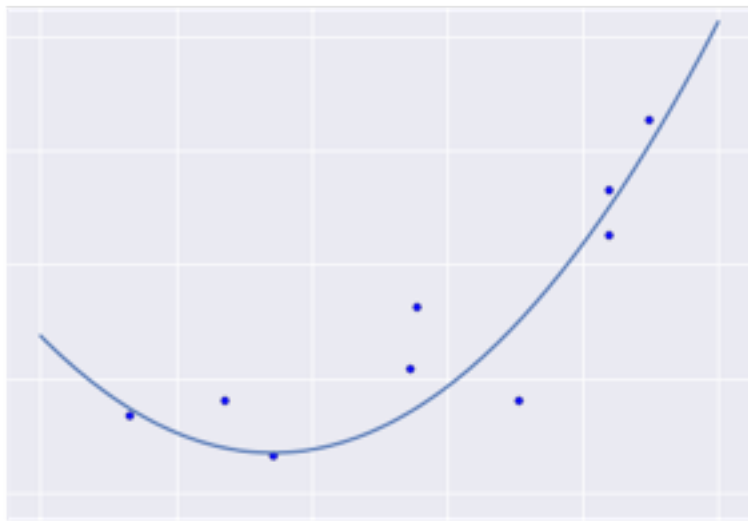
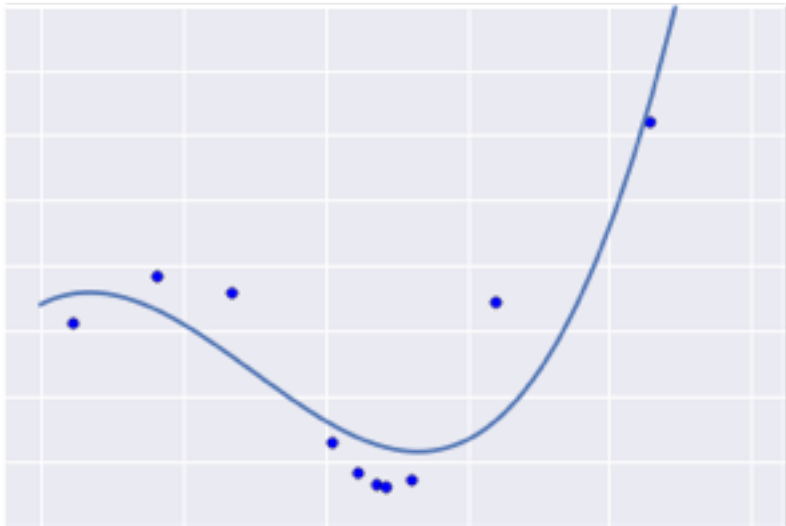
I.

I. REGRESSION RECAP

$$y = \alpha + \beta_1 x_1 + \dots + \beta_n x_n + \varepsilon$$



$$y = \alpha + \beta_1 x_1 + \beta_1 x_1^2 + \dots + \varepsilon$$



OLS:

$$\min (\|y - x\beta\|^2)$$

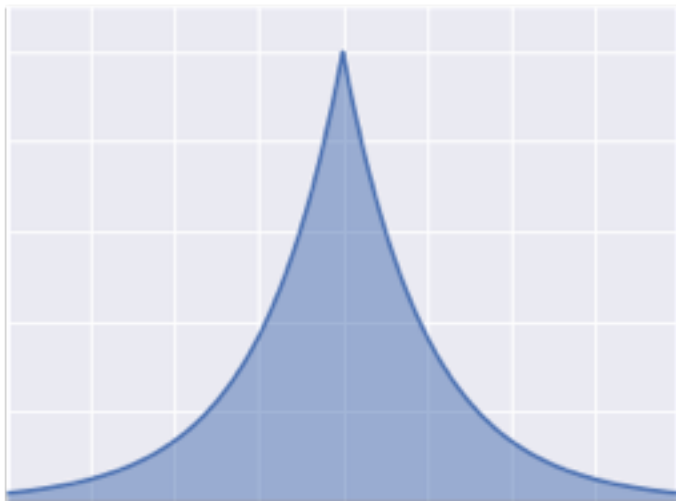
Lasso (L1):

$$\min (\|y - x\beta\|^2 + \lambda \|\beta\|)$$

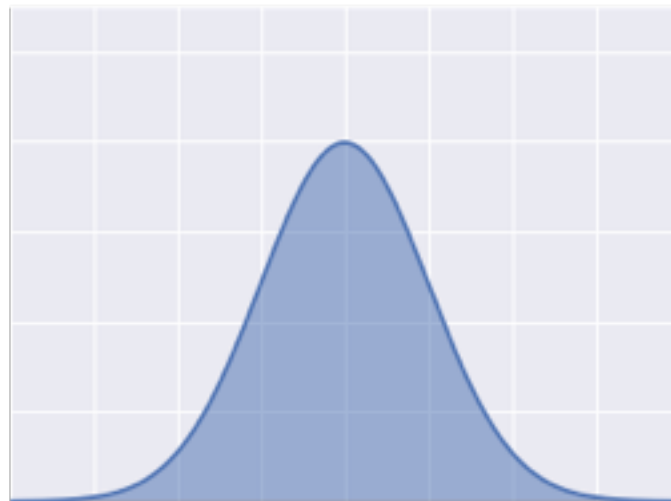
Ridge (L2):

$$\min (\|y - x\beta\|^2 + \lambda \|\beta\|^2)$$

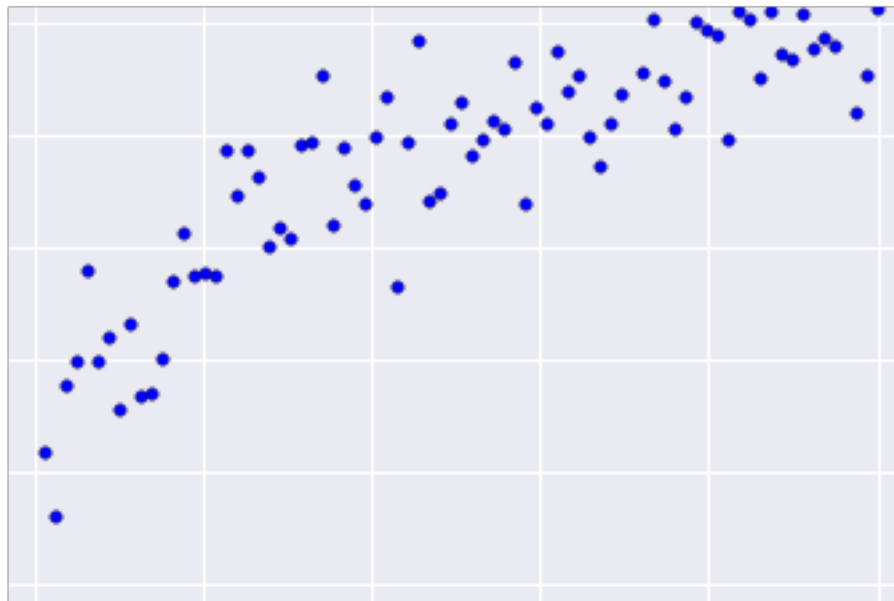
Laplace distribution



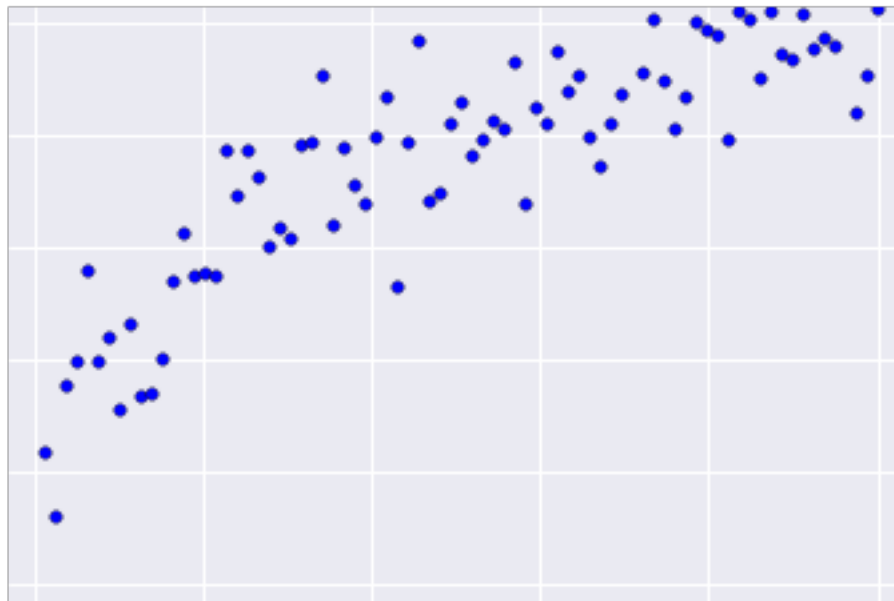
Gaussian distribution



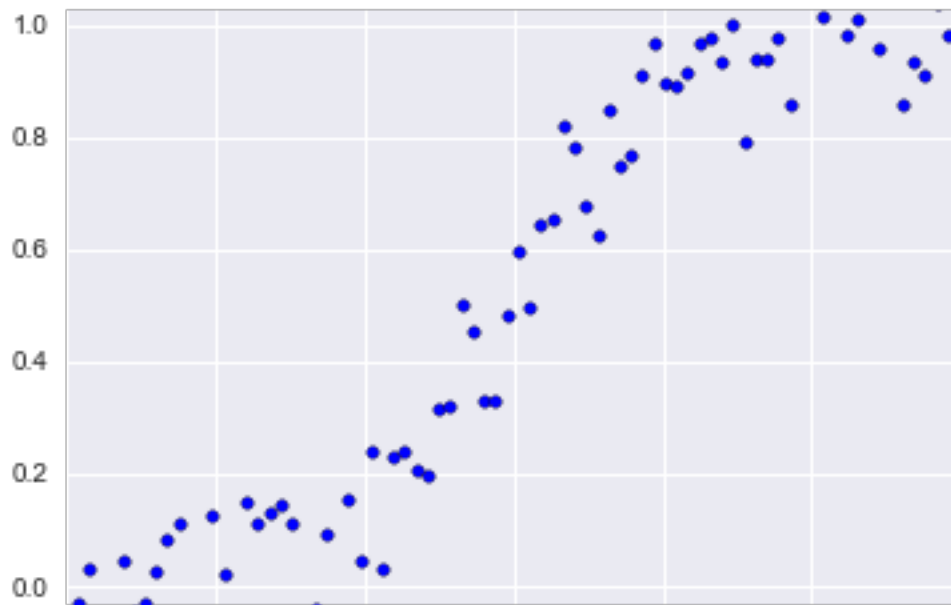
$$y^2 = \alpha + \beta_1 x_1 + \dots + \beta_n x_n + \varepsilon$$



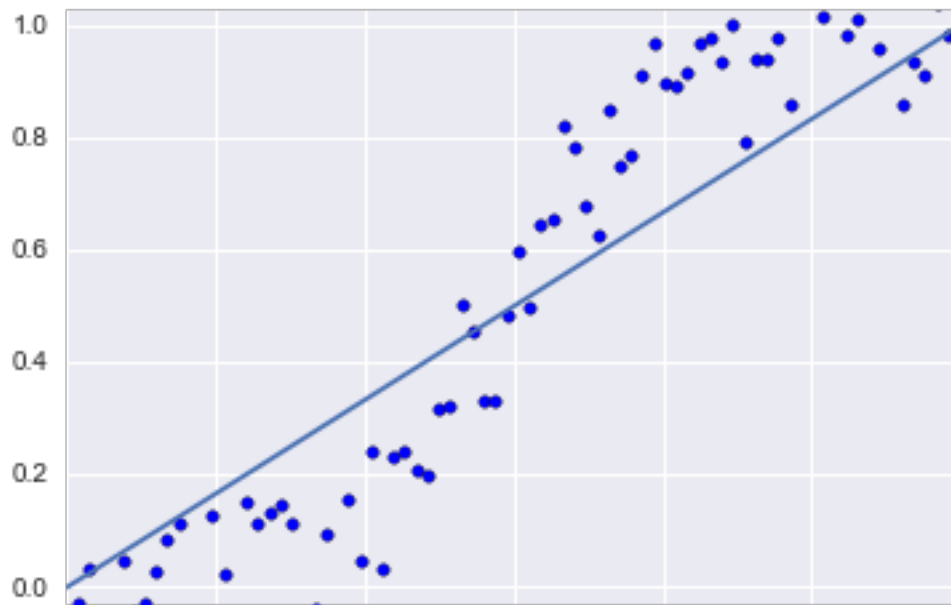
$$\log y = \alpha + \beta_1 x_1 + \dots + \beta_n x_n + \varepsilon$$



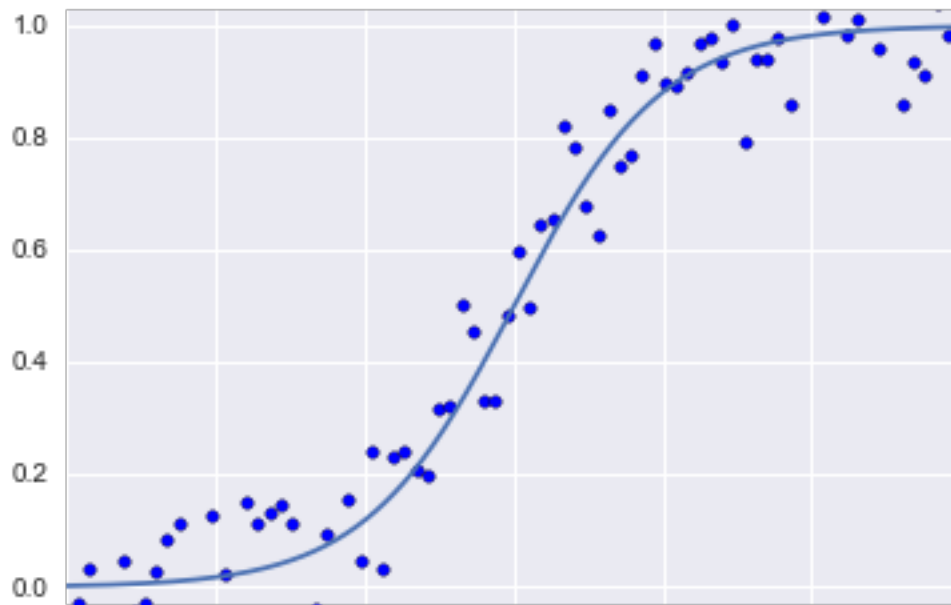
Could we also use regression to estimate probabilities?



Could we also use regression to estimate probabilities?



Could we also use regression to estimate probabilities?



II. LOGISTIC REGRESSION

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

Q: What is logistic regression?

Q: What is logistic regression?

A: A generalization of linear regression to classification problems.

*In **linear regression**, features predict a continuous outcome variable.*

*In **linear regression**, features predict a continuous outcome variable.*

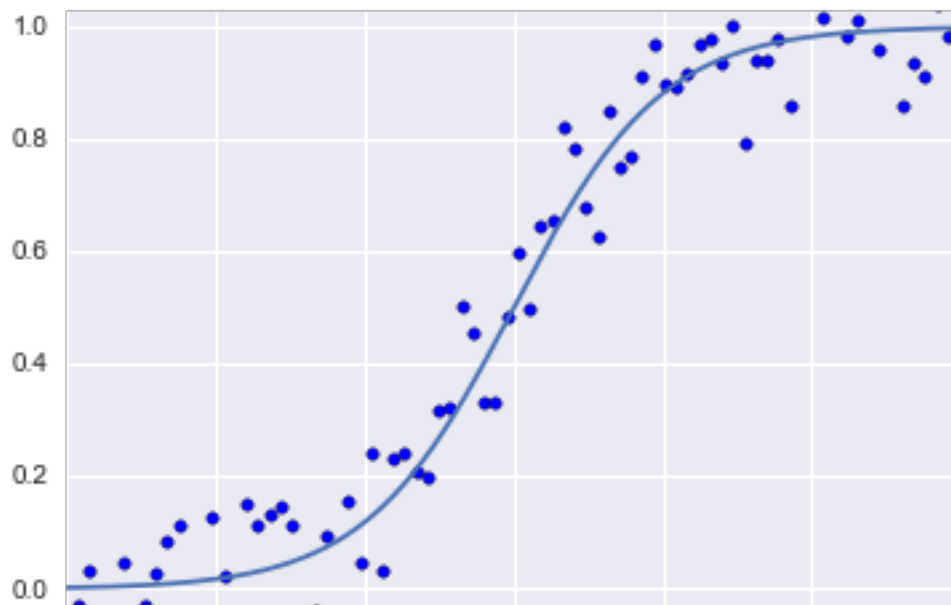
*In **logistic regression**, features predict probabilities of (binary) class membership.*

*In **linear regression**, features predict a continuous outcome variable.*

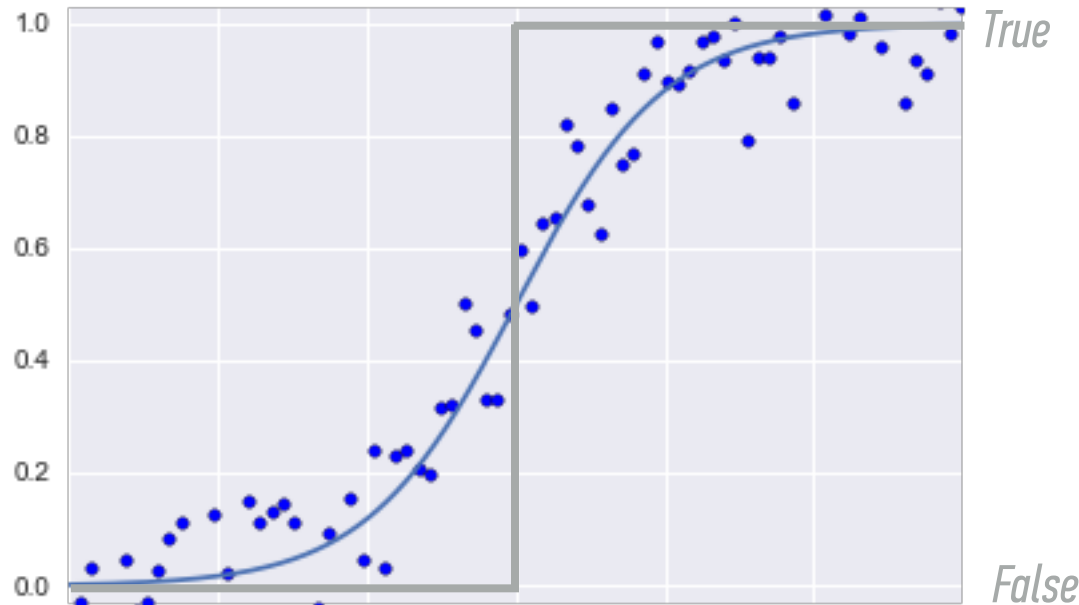
*In **logistic regression**, features predict probabilities of (binary) class membership.*

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

Logistic regression gives us predicted probabilities



Logistic regression gives us predicted probabilities, which then could be ‘snapped’ to class labels



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 first difference is in the outcome variable.

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

The first difference is in the outcome variable.

The second difference is in the error term.

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 $[0, 1]$.

So the first step in extending the linear regression model is to map the outcome variable 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 $[0, 1]$.

So the first step in extending the linear regression model is to map the outcome variable into the unit interval.

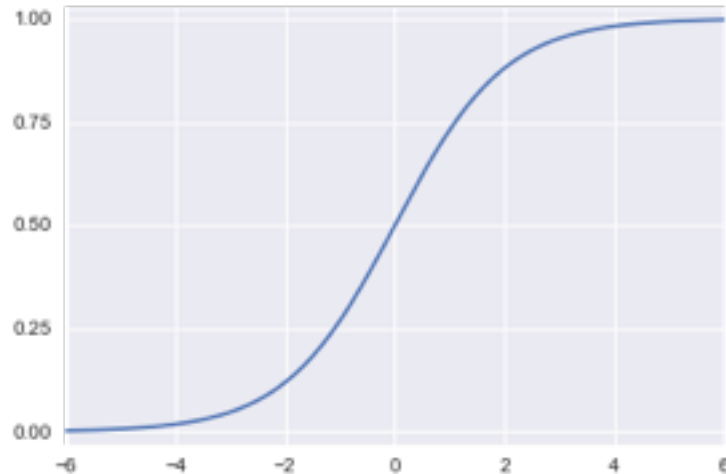
Q: How do we do this?

A: By using a logistic or sigmoid function:

$$\text{sigmoid}(x) = \frac{1}{1 + e^{-x}}$$

A: By using a logistic or sigmoid function:

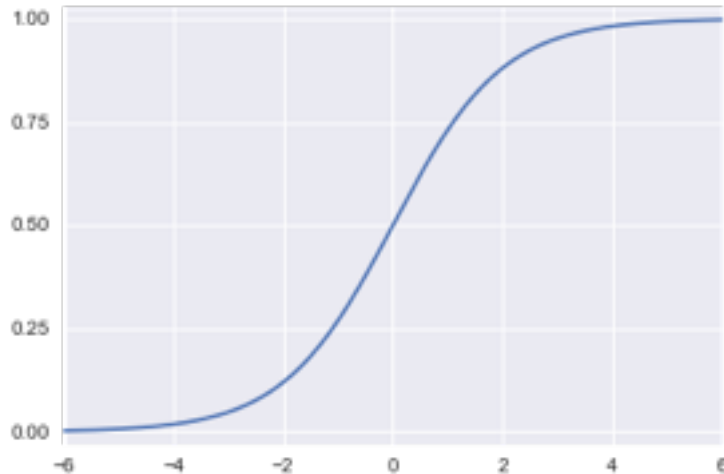
$$\text{sigmoid}(x) = \frac{1}{1 + e^{-x}}$$



*We've already seen
what it looks like*

A: By using a logistic or sigmoid function:

$$\text{sigmoid}(x) = \frac{1}{1 + e^{-x}}$$



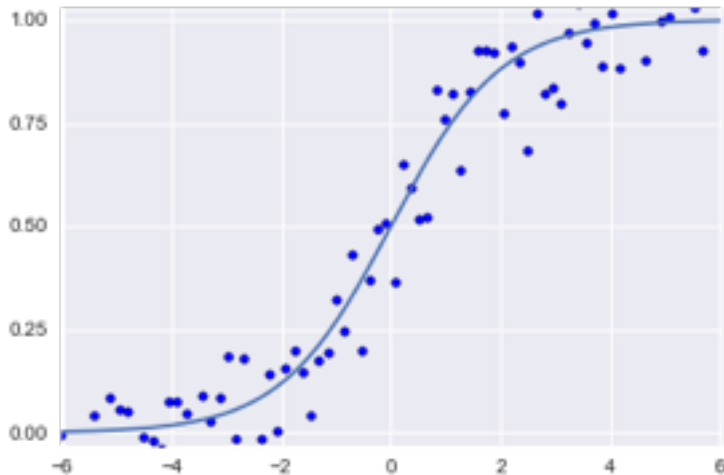
*We've already seen
what it looks like*

NOTE

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

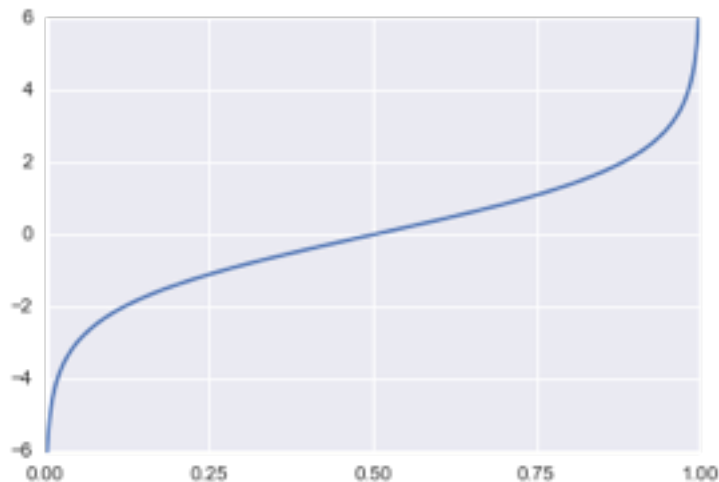
This is a nonlinear
transformation!

$$y = \frac{1}{1 + e^{-(\alpha + \beta_1 x_1 + \dots + \beta_n x_n)}} + \varepsilon$$



*The inverse of the logistic function is the **logit** function,*

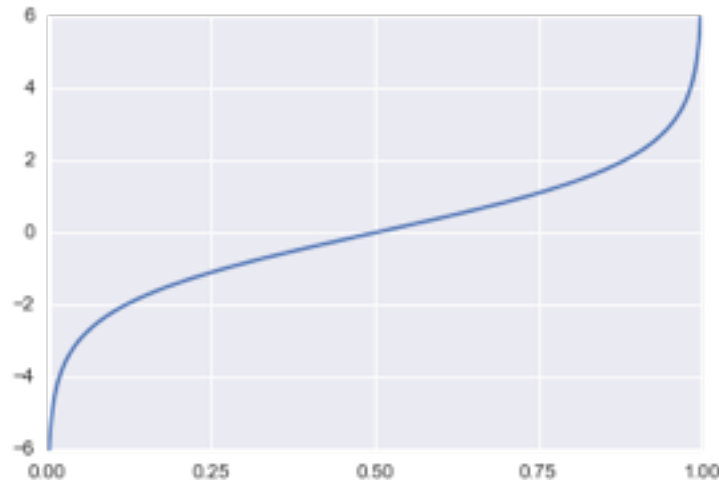
$$\text{logit}(p) = \log \frac{p}{1-p}$$



*The inverse of the logistic function is the **logit** function,*

$$\text{logit}(p) = \log \frac{p}{1-p}$$

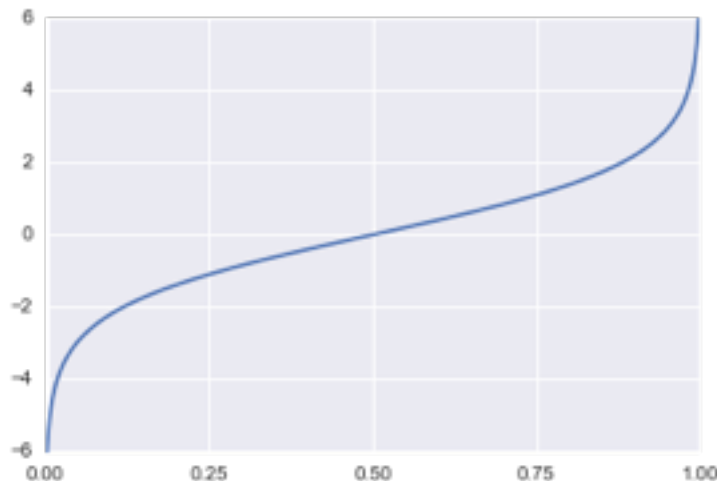
*The logit function is also called the **log-odds** function.*



*The inverse of the logistic function is the **logit** function,*

$$\text{logit}(p) = \log \frac{p}{1-p}$$

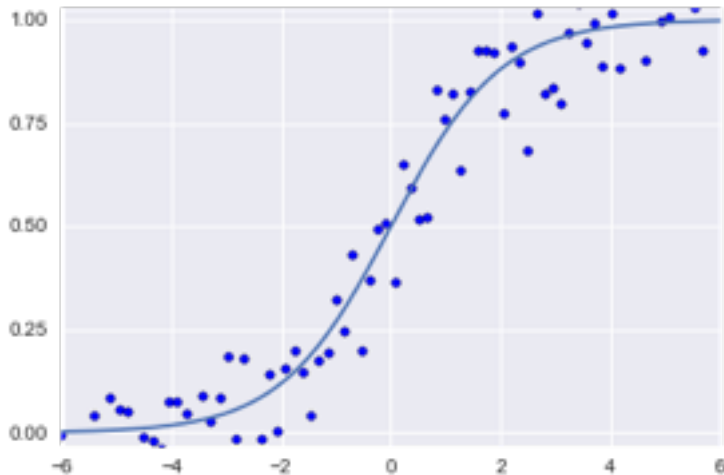
*The logit function is also called the **log-odds** function.*

**NOTE**

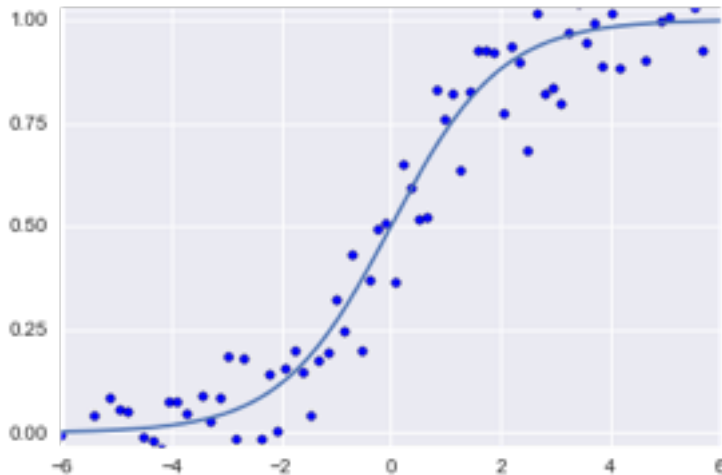
This name hints at its usefulness in interpreting our results.

We will see why shortly.

$$y = \frac{1}{1 + e^{-(\alpha + \beta_1 x_1 + \dots + \beta_n x_n)}} + \varepsilon$$



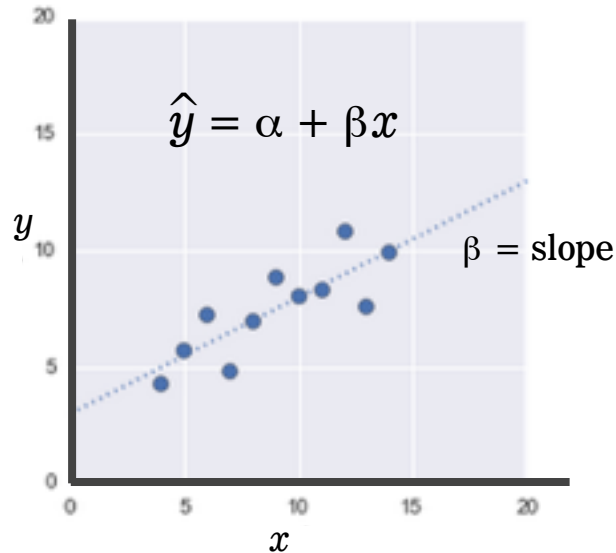
$$\text{logit } y = \alpha + \beta_1 x_1 + \dots + \beta_n x_n + \varepsilon'$$



III. INTERPRETING RESULTS

*In **linear regression**, the parameter β represents the change in the response variable for a unit change in the **X**.*

*In **linear regression**, the parameter β represents the change in the response variable for a unit change in the **X**.*



Example

With each additional year work experience, your salary will grow with β dollars

*In **linear regression**, the parameter β represents the change in the response variable for a unit change in the **X**.*

*In **logistic regression**, β represents the change in the logit function for a unit change in the **X**.*

*In **linear regression**, the parameter β represents the change in the response variable for a unit change in the **X**.*

*In **logistic regression**, β represents the change in the logit function for a unit change in the **X**.*

Interpreting this change in the logit function requires another definition first.

The odds of an event are given by the ratio of the probability of the event by its complement:

$$\text{odds}(p) = \frac{p}{1 - p}$$

The odds of an event are given by the ratio of the probability of the event by its complement:

$$\text{odds}(p) = \frac{p}{1 - p}$$

The odds tell you how much more likely an event is to happen, compared to the event not happening

$$\text{odds of } P(A \text{ happens}) = \frac{P(A \text{ happens})}{P(A \text{ doesn't happen})}$$

*In **linear regression**, the parameter β represents the change in the response variable for a unit change in the **X**.*

*In **logistic regression**, β represents the change in the logit function for a unit change in the **X**.*

*In **linear regression**, the parameter β represents the change in the response variable for a unit change in the **X**.*

*In **logistic regression**, β represents the change in the logit function for a unit change in the **X**.*

Example

With each additional year day of training, you will be e^β as likely to succeed

*In **linear regression**, the parameter β represents the change in the response variable for a unit change in the **X**.*

*In **logistic regression**, β represents the change in the logit function for a unit change in the **X**.*

Example (for $\beta = \log 2$)

*With each additional year day of training,
you will be twice as likely to succeed*

Suppose we are interested in mobile purchase behavior.

y = class label denoting purchase/no purchase

x = binary flag if a user's mobile OS is Apple's iOS

Suppose we are interested in mobile purchase behavior.

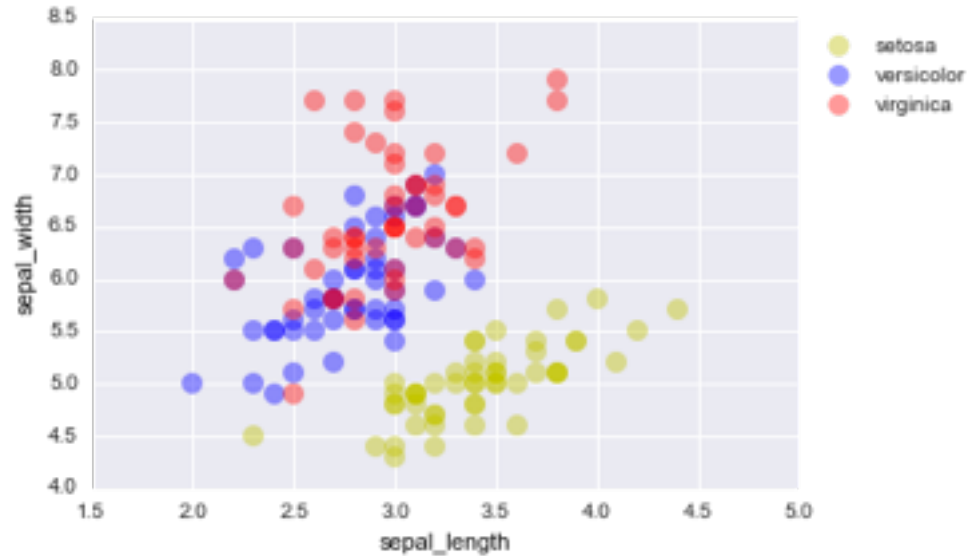
y = class label denoting purchase/no purchase

x = binary flag if a user's mobile OS is Apple's iOS

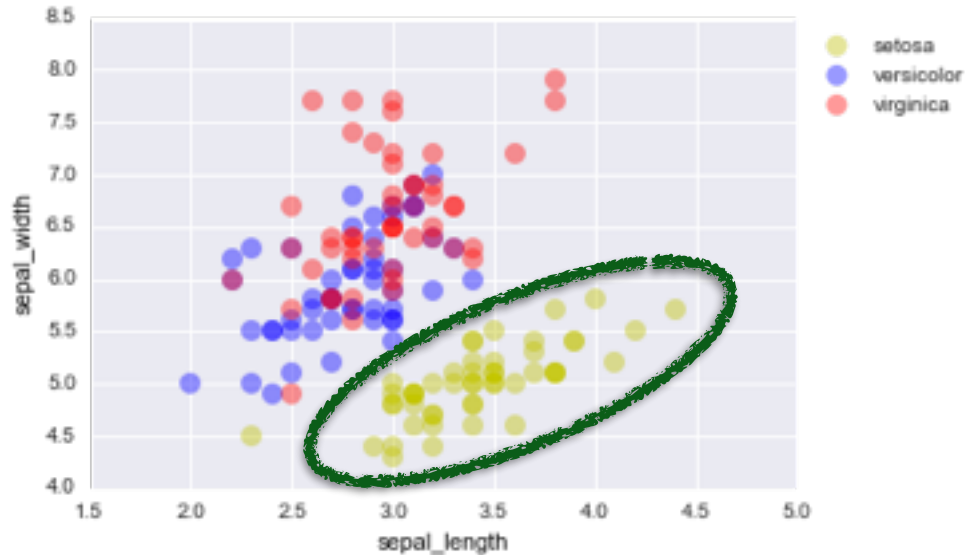
In this case, an odds ratio of 2 (e.g., $\beta = \log 2$) indicates that a purchase is twice as likely for an iOS user as for a non-iOS user.

IV. DECISION BOUNDARIES

Let's have a look at the Iris dataset again

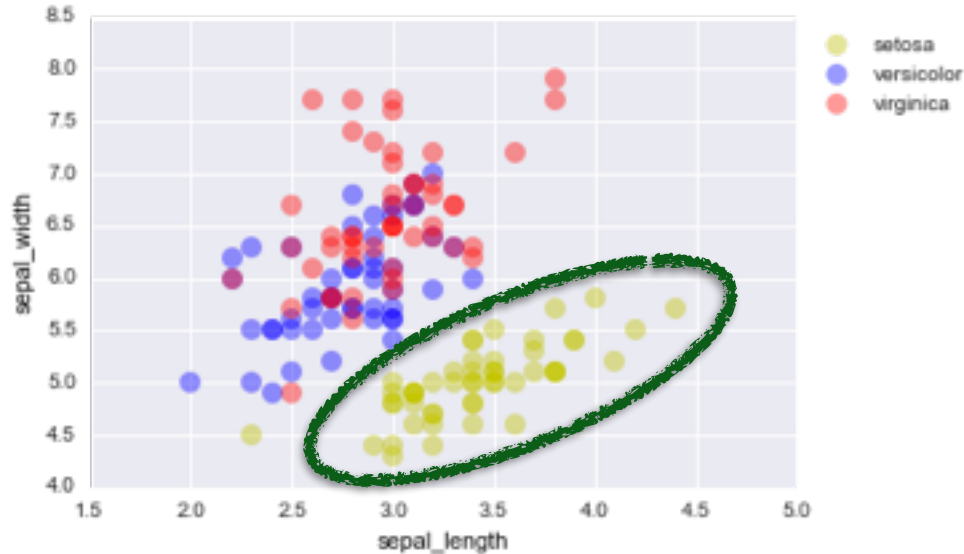


Let's predict, using logistic regression, and only sepal's length and width as features, if a flower is a setosa or not



Let's predict, using logistic regression, and only sepal's length and width as features, if a flower is a setosa or not

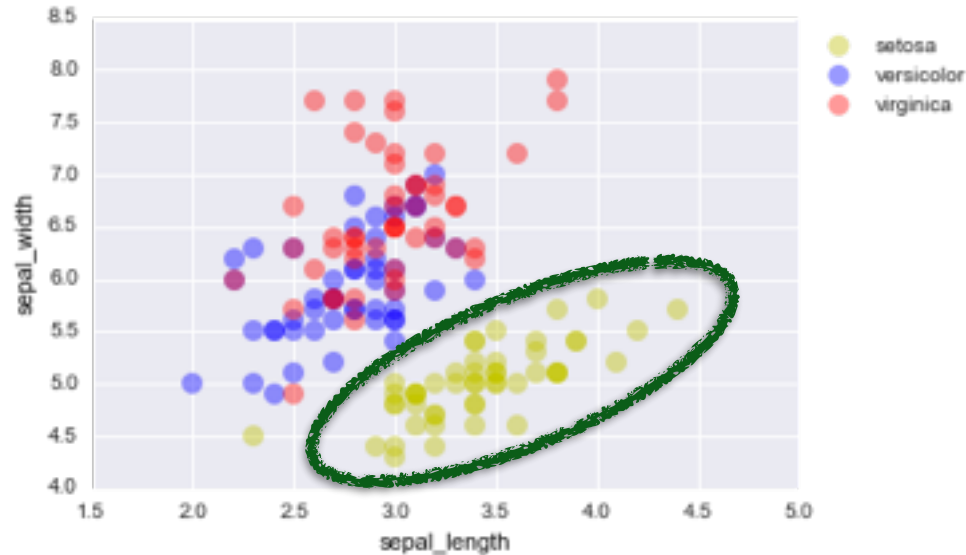
$$\begin{aligned}x_1 &= \text{sepal length} \\x_2 &= \text{sepal width} \\y &= P(\text{setosa})\end{aligned}$$



$$P(\text{setosa}) = \text{sigmoid}(\alpha + \beta_1 x_1 + \beta_2 x_2)$$

$x_1 = \text{sepal length}$

$x_2 = \text{sepal width}$



$$P(\text{setosa}) = \text{sigmoid}(\alpha + \beta_1 x_1 + \beta_2 x_2)$$

If we set our cut-off at 50%, we predict a positive iff $P \geq 0.5$

$$P(\text{setosa}) = \text{sigmoid}(\alpha + \beta_1 x_1 + \beta_2 x_2)$$

If we set our cut-off at 50%, we predict a positive iff $P \geq 0.5$

*The case $P = 0.5$ denotes the **boundary** between our predictions*

$$1/2 = \text{sigmoid}(\alpha + \beta_1 x_1 + \beta_2 x_2)$$

If we set our cut-off at 50%, we predict a positive iff $P \geq 0.5$

*The case $P = 0.5$ denotes the **boundary** between our predictions*

$$0 = \alpha + \beta_1 x_1 + \beta_2 x_2$$

If we set our cut-off at 50%, we predict a positive iff $P \geq 0.5$

*The case $P = 0.5$ denotes the **boundary** between our predictions*

But since $\text{sigmoid}(x) = \frac{1}{1 + e^{-x}}$ it follows that $x = 0$

In general, for a logistic regression model given by

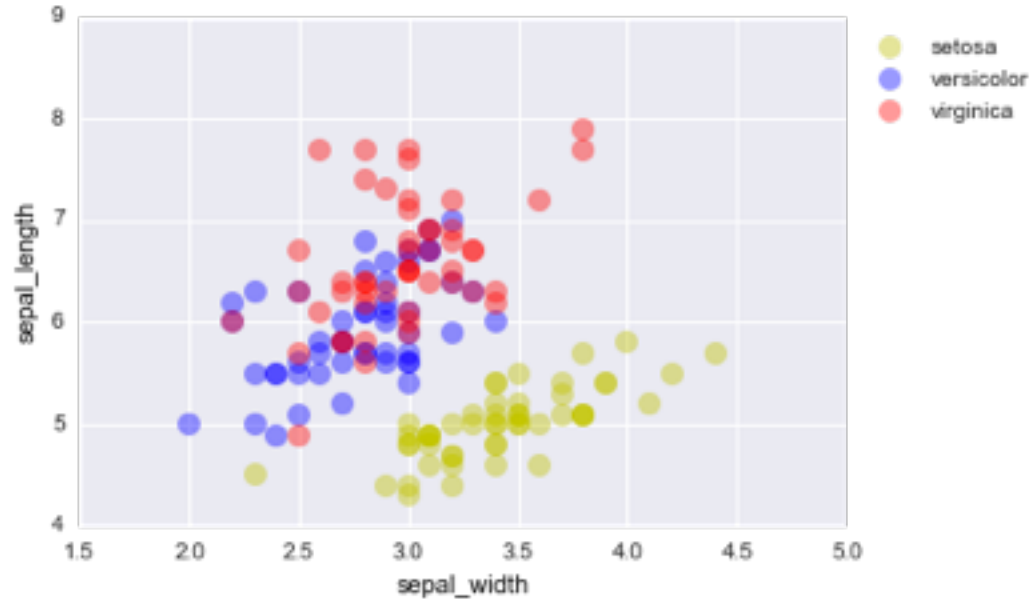
$$\hat{y} = \frac{1}{1 + e^{-(\alpha + \beta_1 x_1 + \dots + \beta_n x_n)}}$$

*its **decision boundary** is given by the equation*

$$\alpha + \beta_1 x_1 + \dots + \beta_n x_n = 0$$

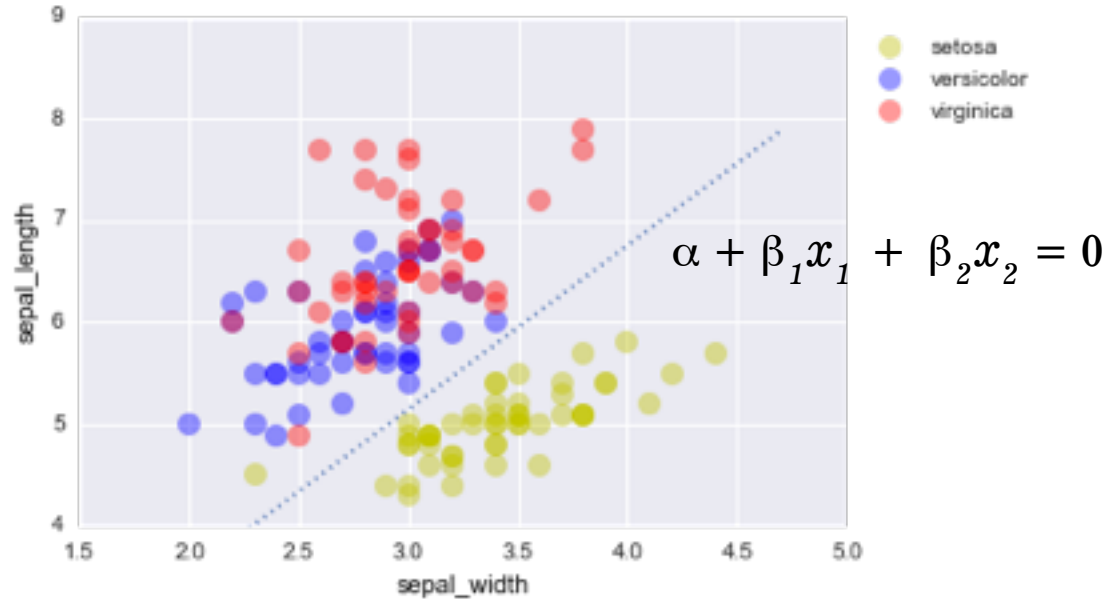
$$P(\text{setosa}) = \text{sigmoid}(\alpha + \beta_1 x_1 + \beta_2 x_2)$$

$x_1 = \text{sepal length}$
 $x_2 = \text{sepal width}$
 $y = P(\text{setosa})$



$$P(\text{setosa}) = \text{sigmoid}(\alpha + \beta_1 x_1 + \beta_2 x_2)$$

$x_1 = \text{sepal length}$
 $x_2 = \text{sepal width}$
 $y = P(\text{setosa})$



V. EVALUATING CLASSIFIERS

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

Q: How do we evaluate a classification model?

Q: How do we evaluate a classification model?

A. Accuracy

of times we make the correct classification / # classifications

Q: How do we evaluate a classification model?

A. Accuracy

of times we make the correct classification / # classifications

When is this a bad a metric?

Q: How do we evaluate a classification model?

A. Accuracy

of times we make the correct classification / # classifications

When is this a bad a metric?

A: When predicting rare or very likely events

predictions

Yes

No

		<i>predictions</i>	
		<i>Yes</i>	<i>No</i>
<i>observations</i>	<i>Yes</i>		
	<i>No</i>		

		<i>predictions</i>	
		<i>Yes</i>	<i>No</i>
<i>observations</i>	<i>Yes</i>	true positive	
	<i>No</i>		

		<i>predictions</i>	
		<i>Yes</i>	<i>No</i>
<i>observations</i>	<i>Yes</i>	true positive	
	<i>No</i>		true negative

		<i>predictions</i>	
		<i>Yes</i>	<i>No</i>
<i>observations</i>	<i>Yes</i>	true positive	false negative
	<i>No</i>		true negative

		<i>predictions</i>	
		<i>Yes</i>	<i>No</i>
<i>observations</i>	<i>Yes</i>	true positive	false negative
	<i>No</i>	false positive	true negative

$$\text{Accuracy} = (TP + TN) / \text{all}$$

observ.	<i>predictions</i>	
	<i>Yes</i>	<i>No</i>
<i>Yes</i>	TP	FN
<i>No</i>	FP	TN

$$\text{Accuracy} = (TP + TN) / \text{all}$$

$$\text{Precision} = TP / (TP + FP)$$

		<i>predictions</i>	
		<i>Yes</i>	<i>No</i>
<i>observ.</i>	<i>Yes</i>	TP	FN
	<i>No</i>	FP	TN

“Of all the cases I highlighted, how often was I right?”

Accuracy = $(TP + TN) / \text{all}$

Precision = $TP / (TP + FP)$

Recall = $TP / (TP + FN)$
aka *hit rate* or *sensitivity*

		<i>predictions</i>	
		<i>Yes</i>	<i>No</i>
<i>observ.</i>	<i>Yes</i>	TP	FN
	<i>No</i>	FP	TN

“Of all the cases to be highlighted, how many did I hit?”

Accuracy = $(TP + TN) / all$

Precision = % correct predictions of all positive predictions

Recall = % correct predictions of all positive cases

Accuracy = (TP + TN) / all

Precision = % correct predictions of all positive predictions

Recall = % correct predictions of all positive cases

F1 score = $2 \frac{P \times R}{P + R}$

Q: When do you want a high recall model?

Q: When do you want a high recall model?

A. Cost of false positive is low, cost of false negative is high

i.e. Predict who should receive a new cheap drug with low side effects

We want to capture as many people as we can

Q: When do you want a high precision model?

Q: When do you want a high precision model?

A. Cost of false positive is high

i.e. Predict who should receive an expensive, complicated surgery treatment

We want to make sure we are correct on any predictions

INTRO TO DATA SCIENCE

DICSUSSION