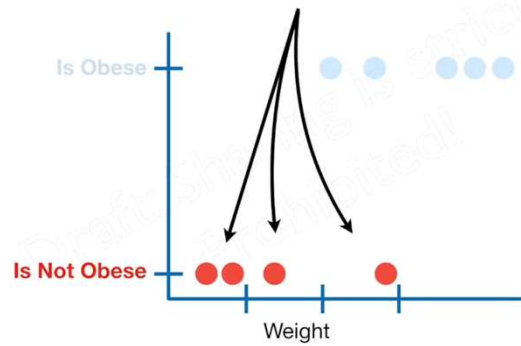
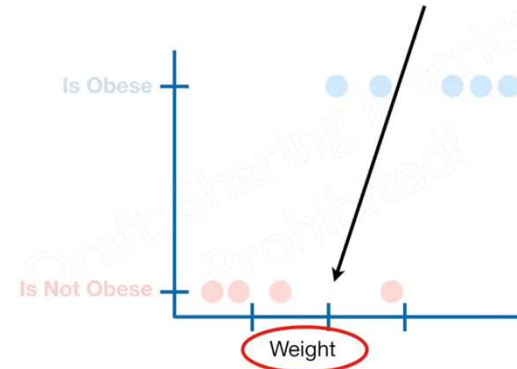


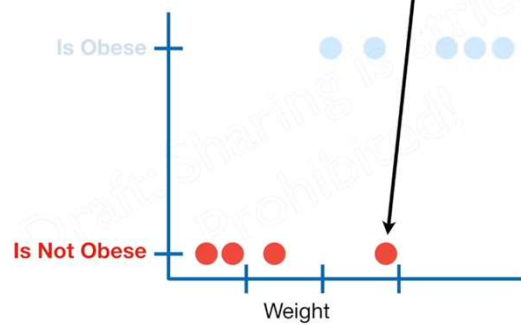
...and the **red dots** represent mice that are **not obese**.



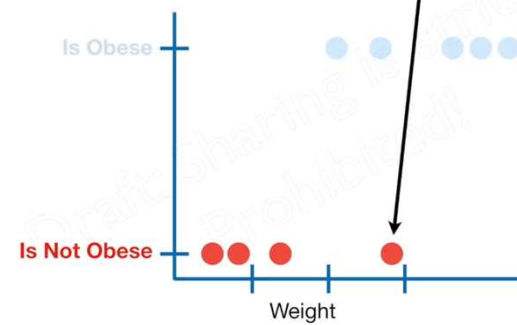
Along the x-axis, we have weight.

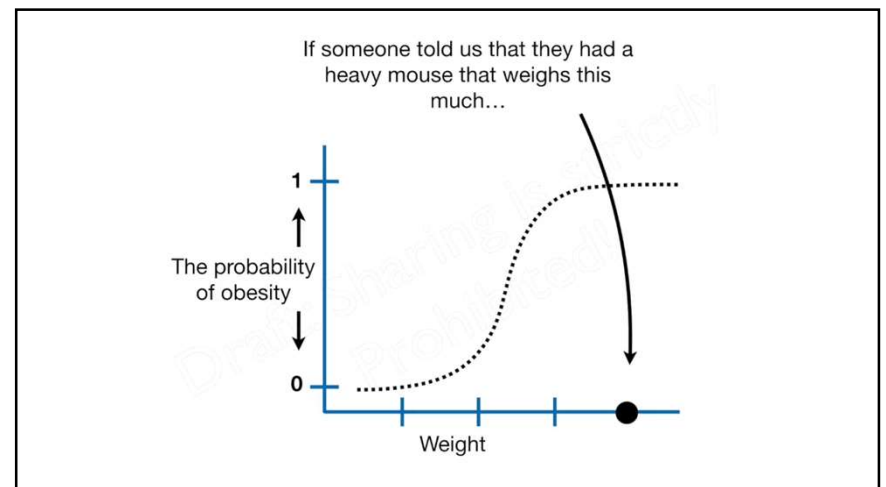
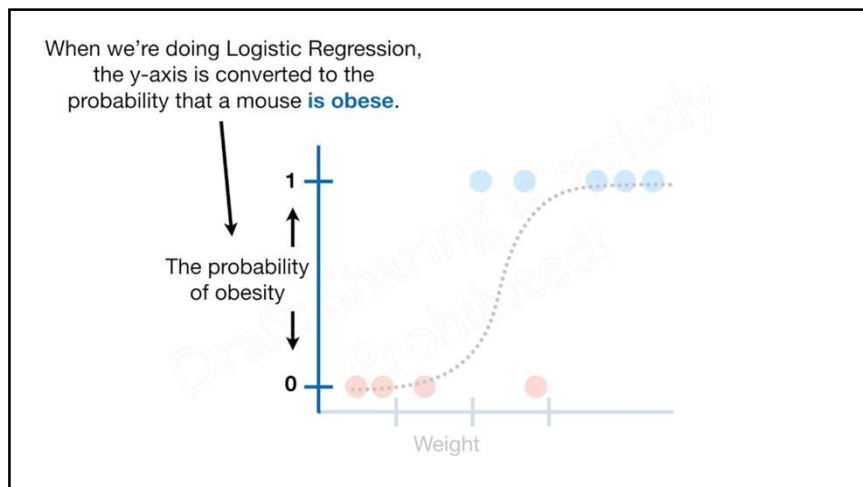
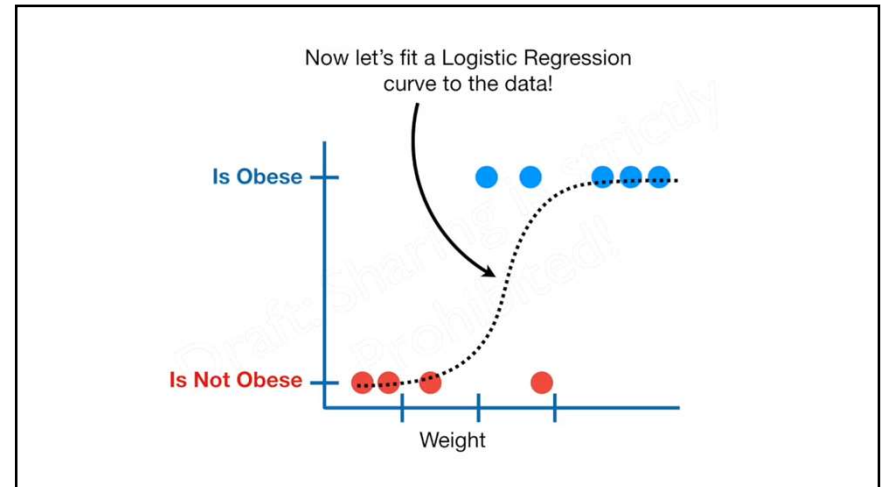
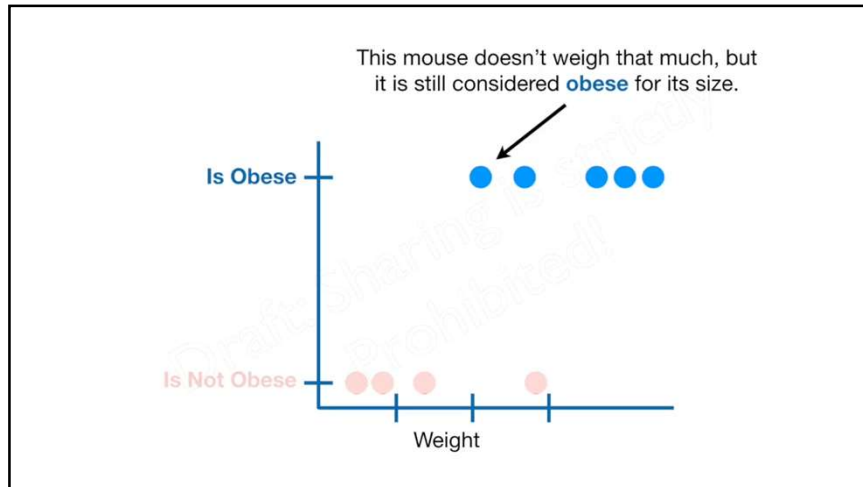


This mouse is **not obese**, even though it weighs a lot.

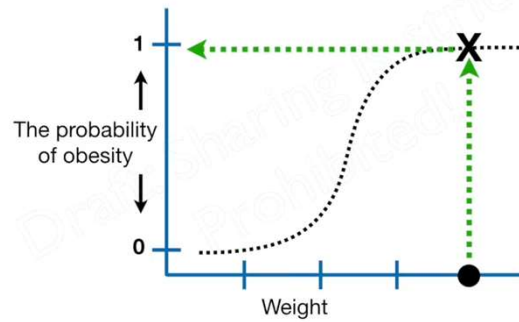


It must be Mighty Mouse and just full of muscles.

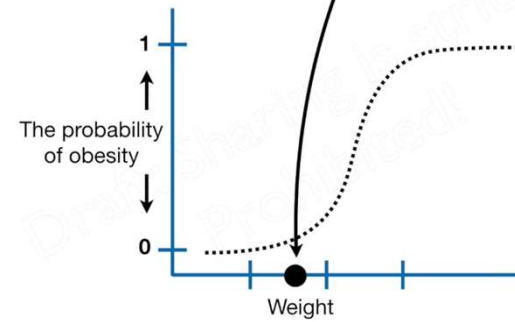




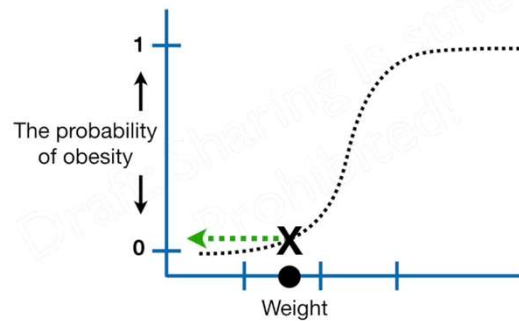
...then the curve would tell us that there is a **high** probability that the mouse **is obese**.



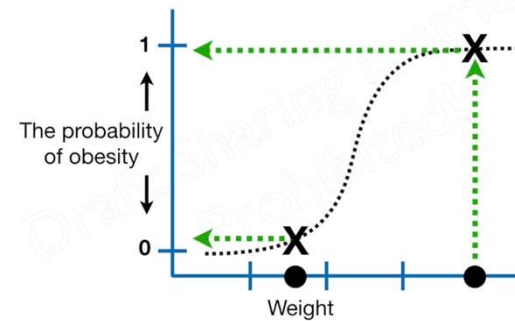
If someone told us that they had a light mouse that weighs this much...



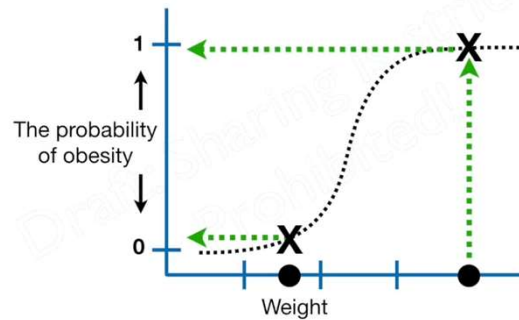
...then the curve would tell us that there is a **low** probability that the mouse **is obese**.



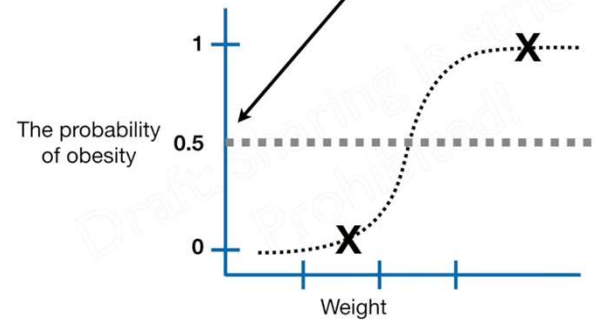
So this Logistic Regression tells us the **probability** that a mouse is **obese** based on its weight.



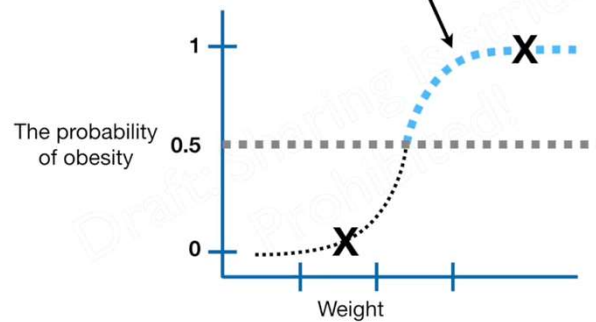
However, if we want to **classify** the mice as **obese** or **not obese**, then we need a way to turn probabilities into classifications.



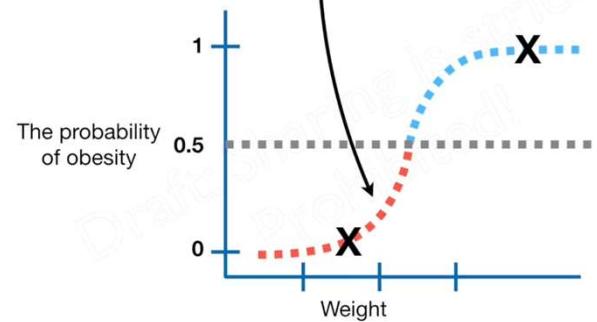
One way to classify mice is to set a threshold at 0.5...



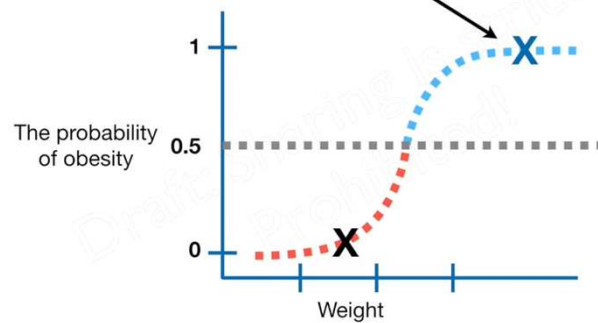
...and classify all mice with a probability of being **obese** > 0.5 as "**obese**"...



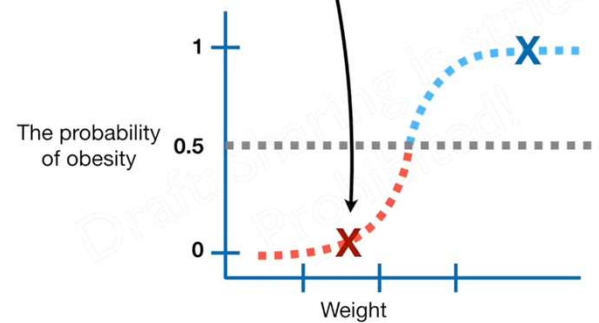
...and classify all mice with a probability of being **obese** ≤ 0.5 as "**not obese**".



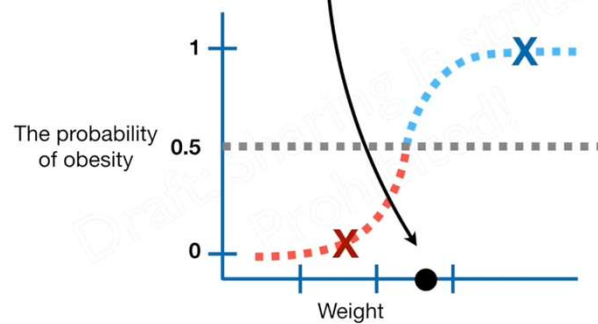
Using **0.5** as the cutoff, we would call this mouse **obese**...



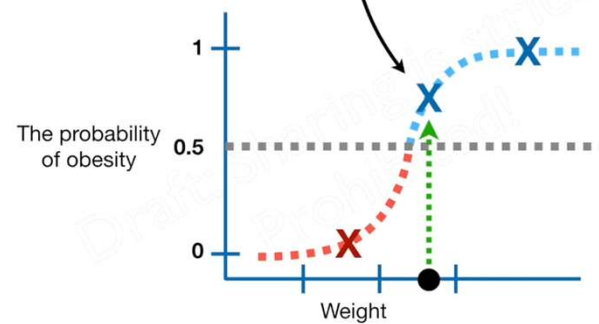
...and this mouse **not obese**.

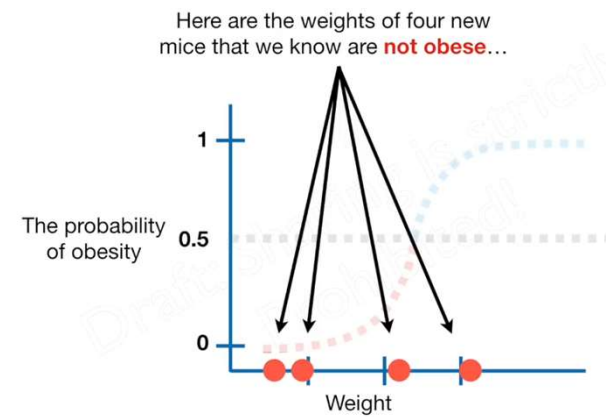
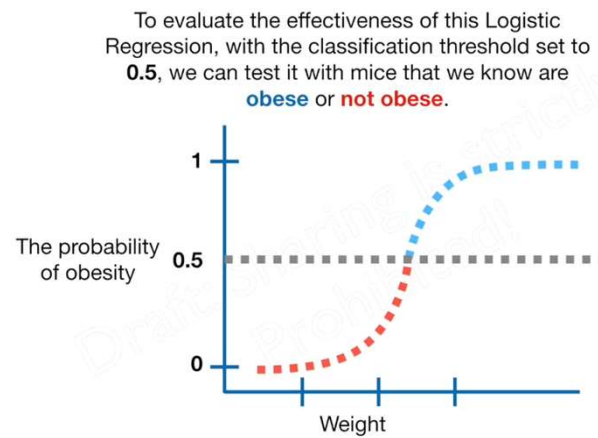
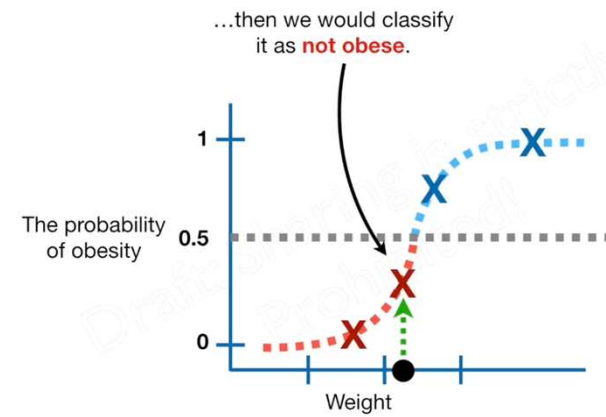
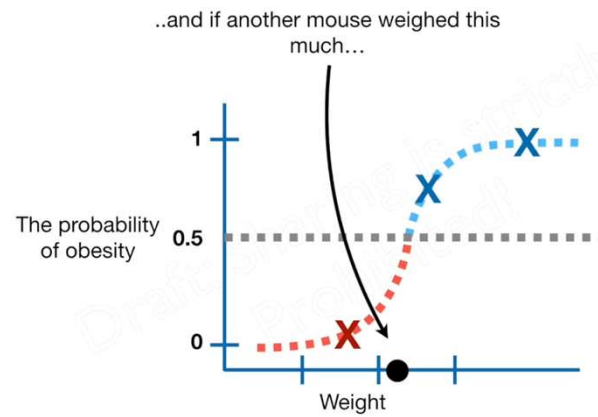


If another mouse weighed this much...

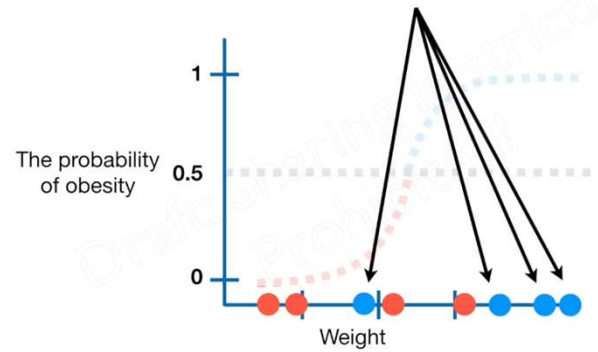


...then we would classify it as **obese**...

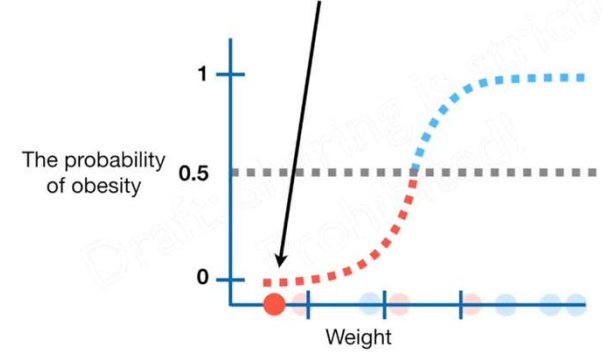




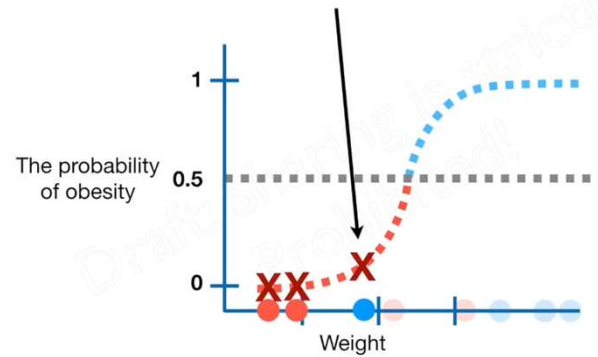
...and here are the weights of four new mice that we know are **obese**.



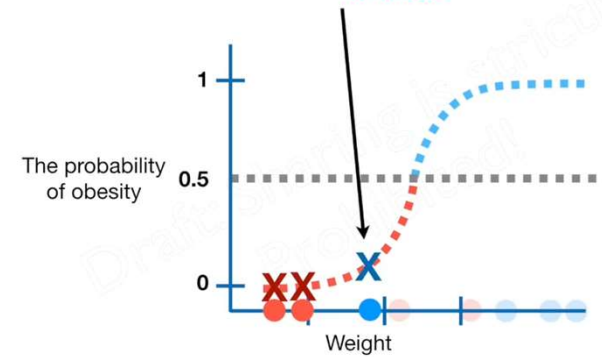
We know that this mouse is **not obese**...

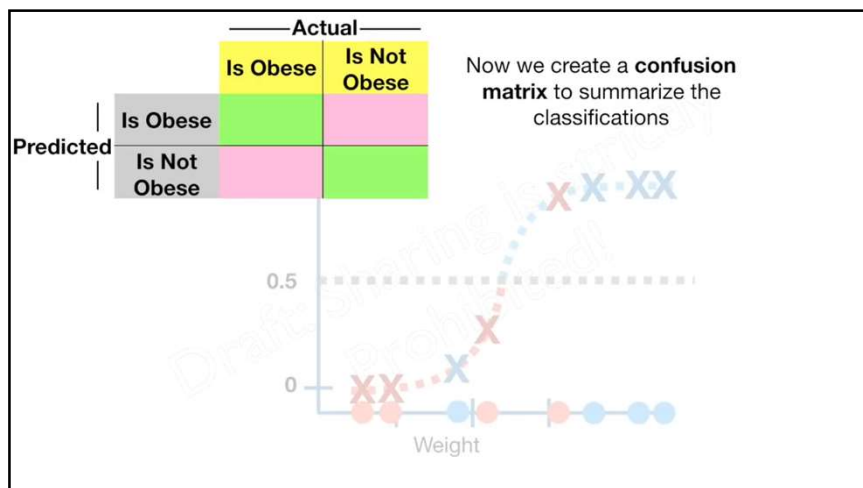
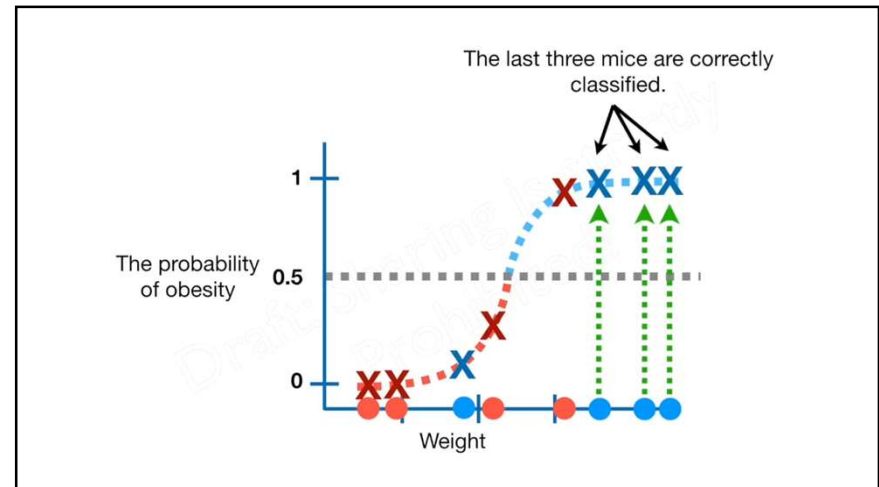
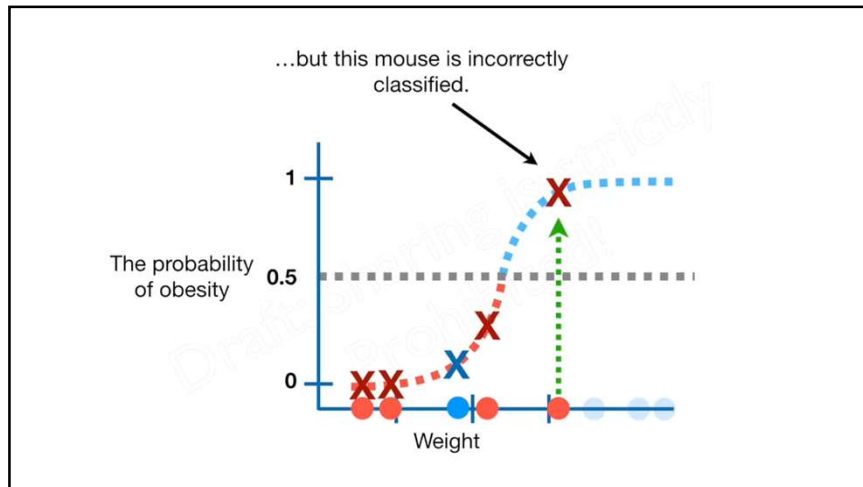


...but this mouse is *incorrectly* classified.

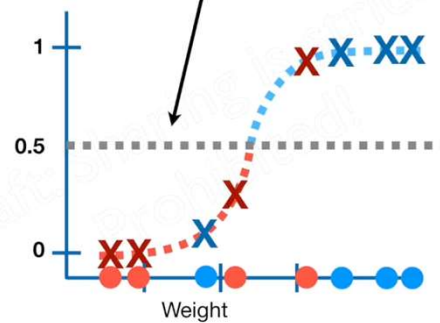


We know that it is **obese**, but it is classified as **not obese**.

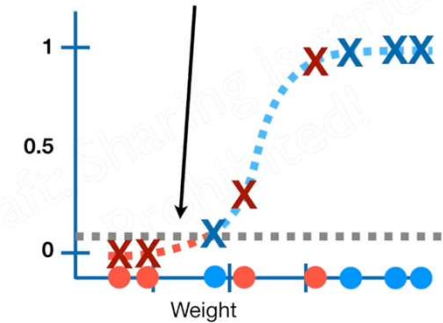




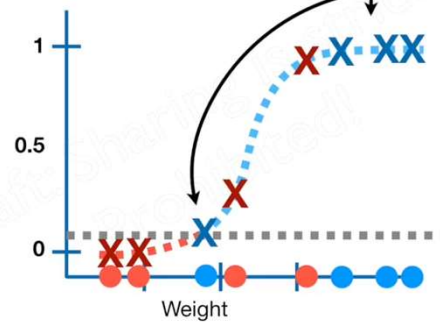
Now let's talk about what happens when we use a different threshold for deciding if a sample is **obese** or **not**.



For example, if it was super important to correctly classify every **obese** sample, we could set the threshold to **0.1**...

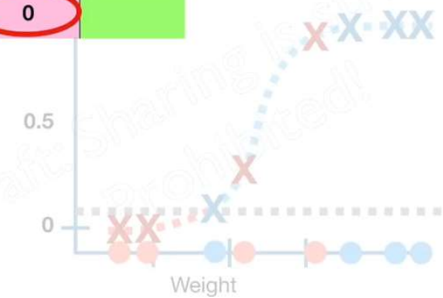


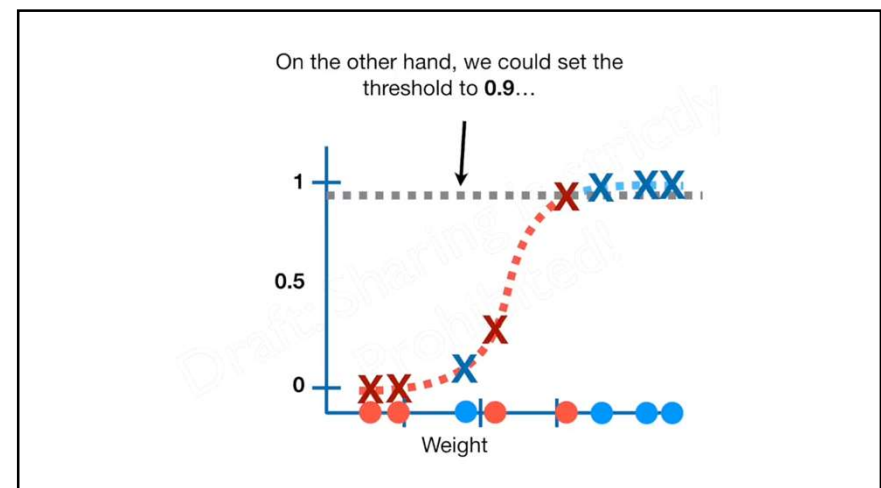
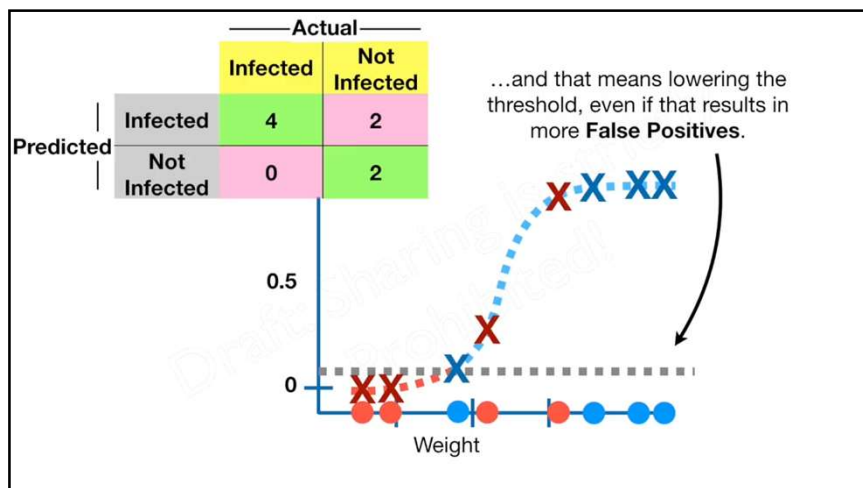
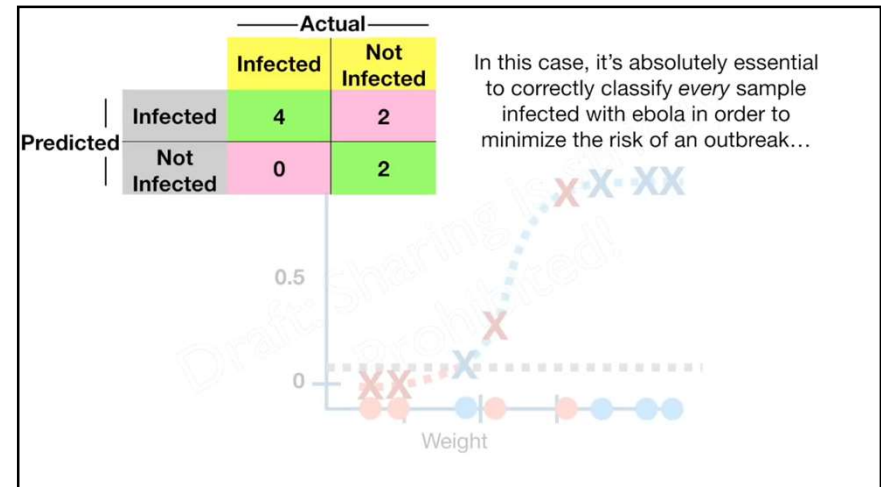
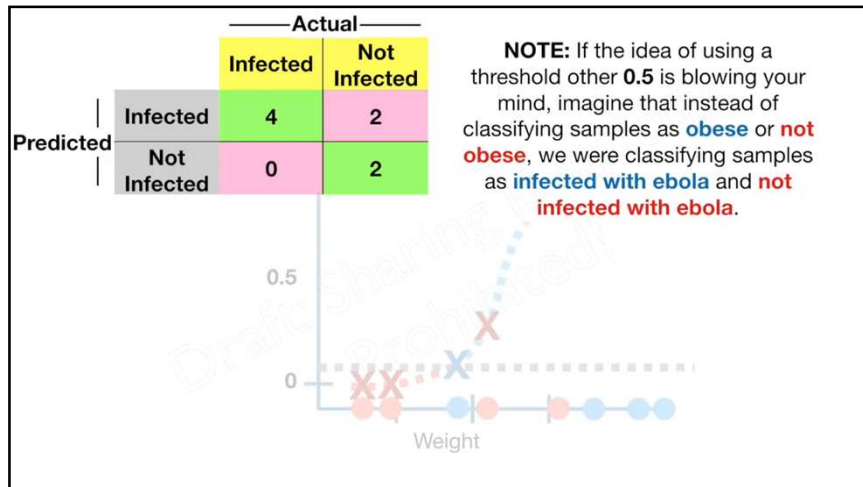
This would result in correct classifications for all 4 **obese** mice...



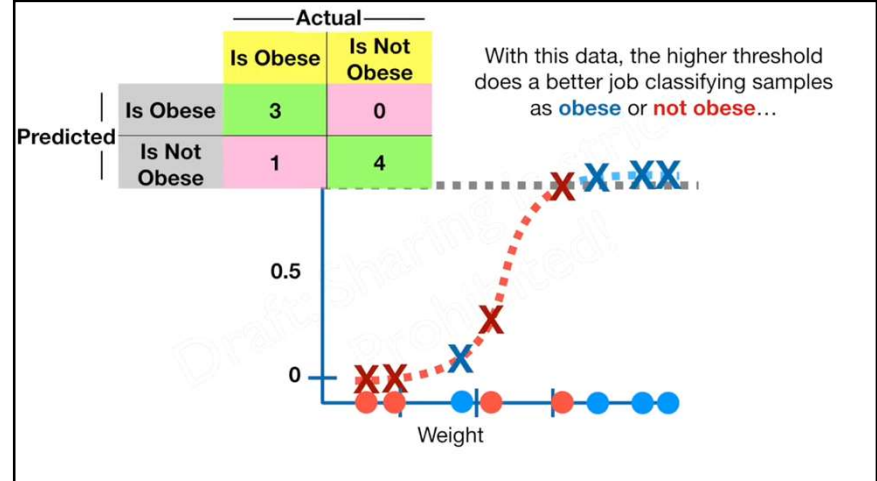
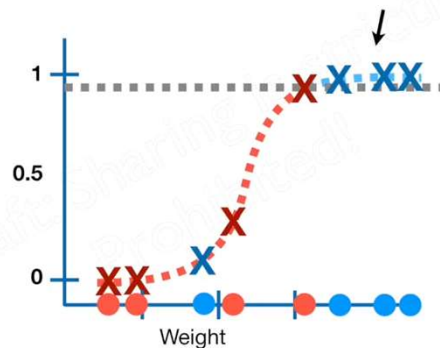
		Actual	
		Is Obese	Is Not Obese
Predicted	Is Obese	4	2
	Is Not Obese	0	

The lower threshold would also reduce the number of **False-Negatives**, because all of the **obese** mice were correctly classified...

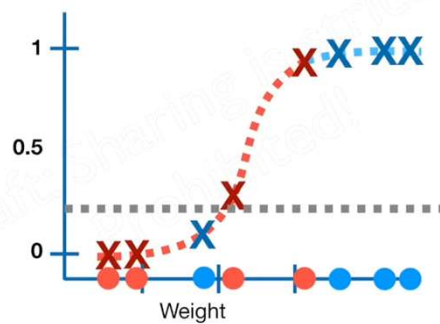




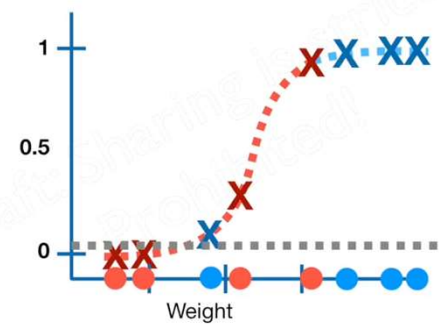
In this case, we would correctly classify the same number of **obese** samples as when the threshold was set to 0.5...



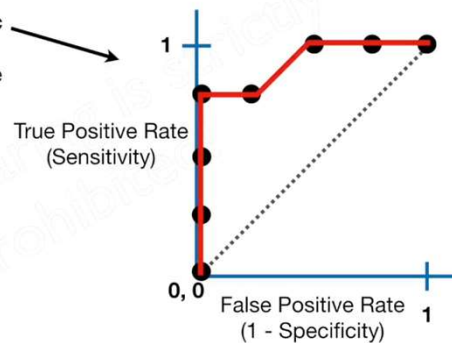
...but the threshold could be set to anything between 0 and 1.



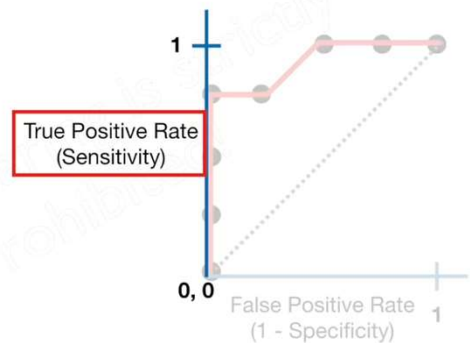
How do we determine which threshold is the best?



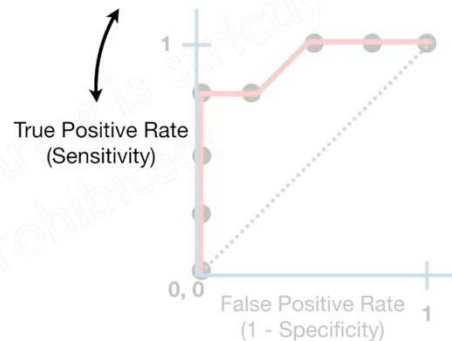
So instead of being overwhelmed with confusion matrices, **Receiver Operator Characteristic (ROC)** graphs provide a simple way to summarize all of the information.



The y-axis shows the **True Positive Rate**, which is the same thing as **Sensitivity**.



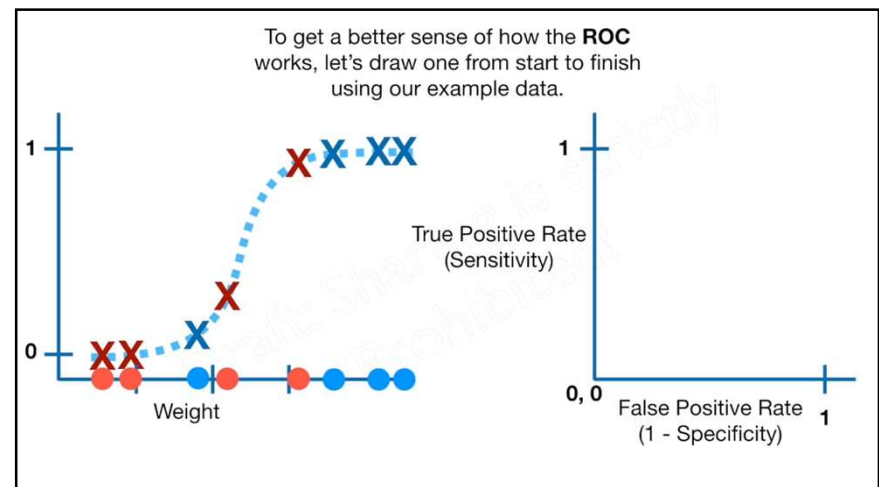
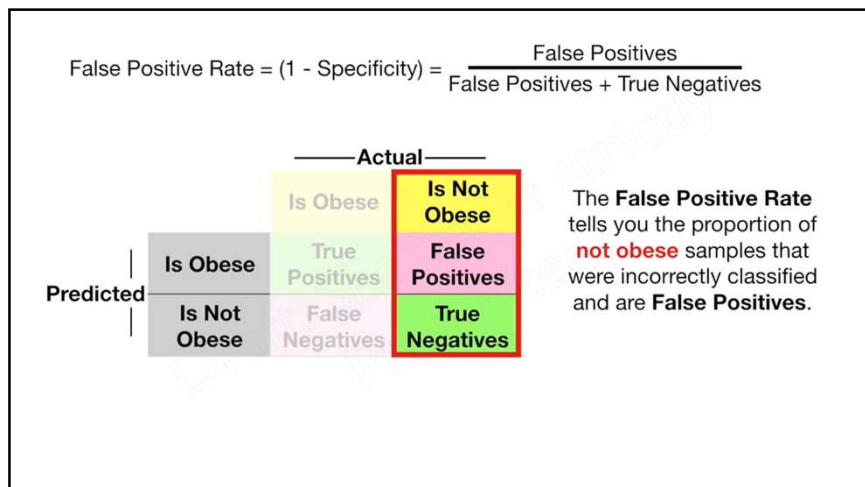
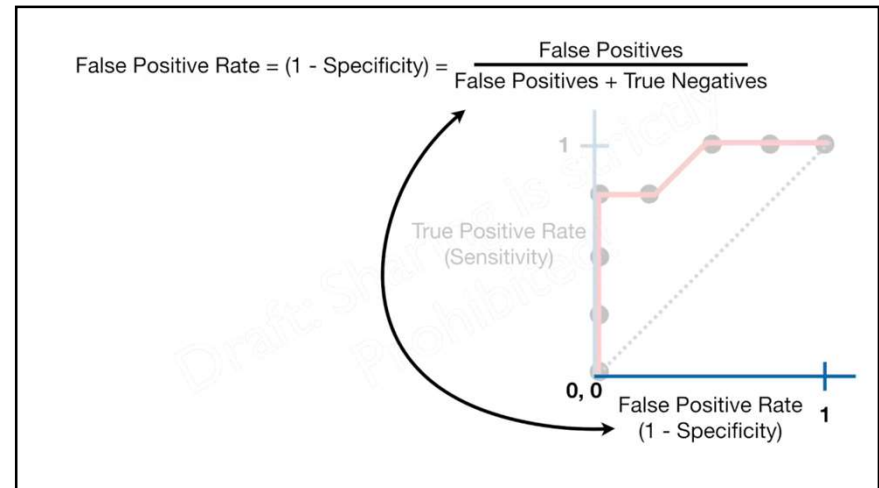
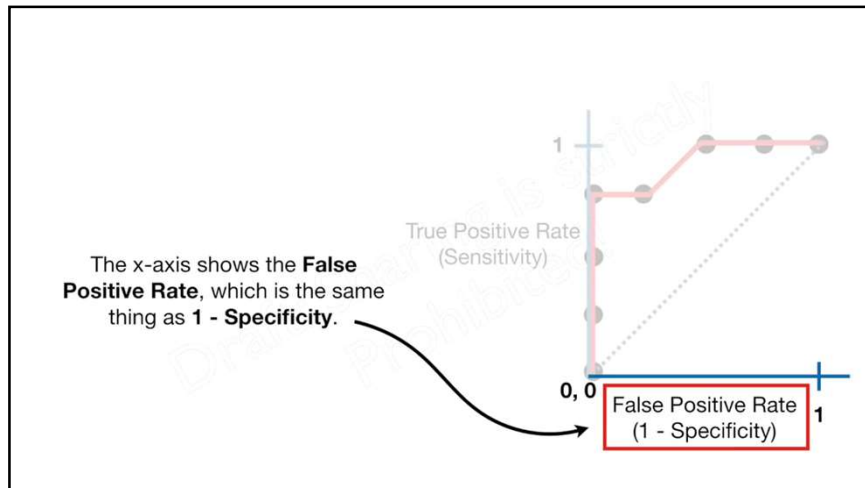
$$\text{True Positive Rate} = \text{Sensitivity} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

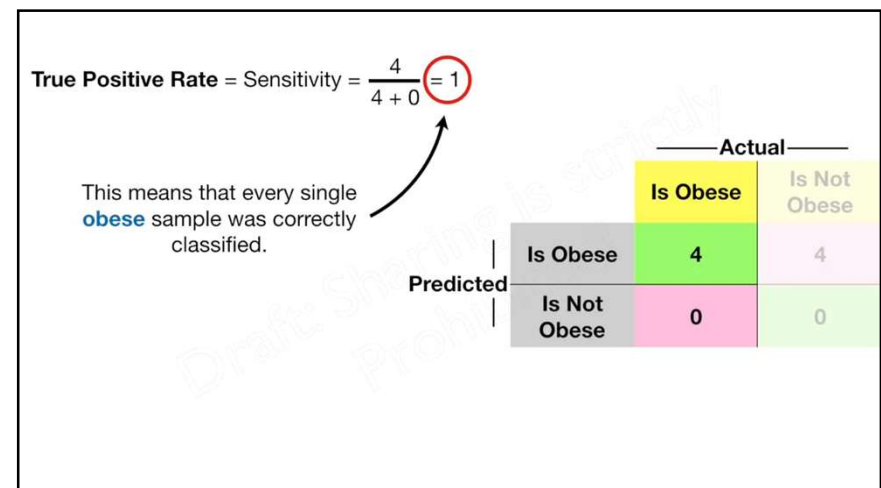
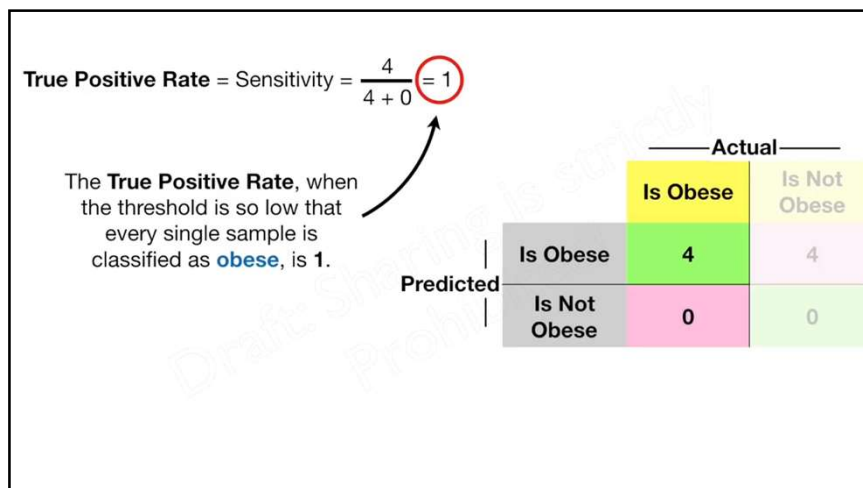
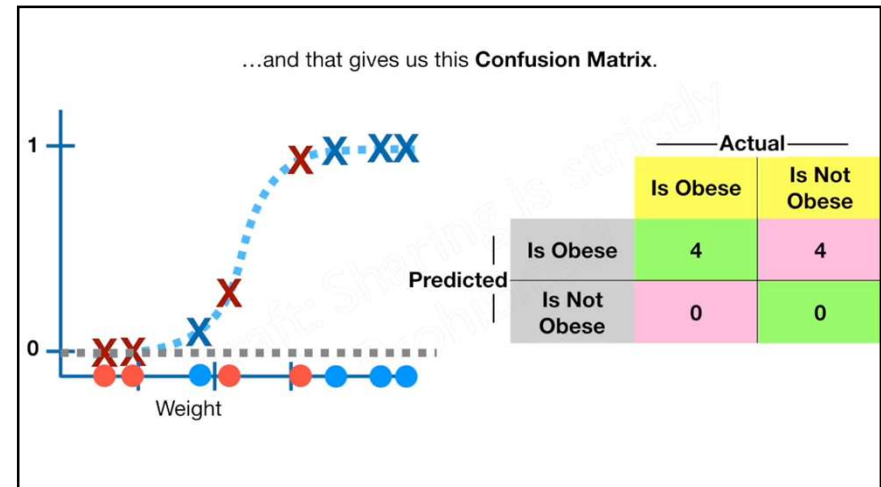
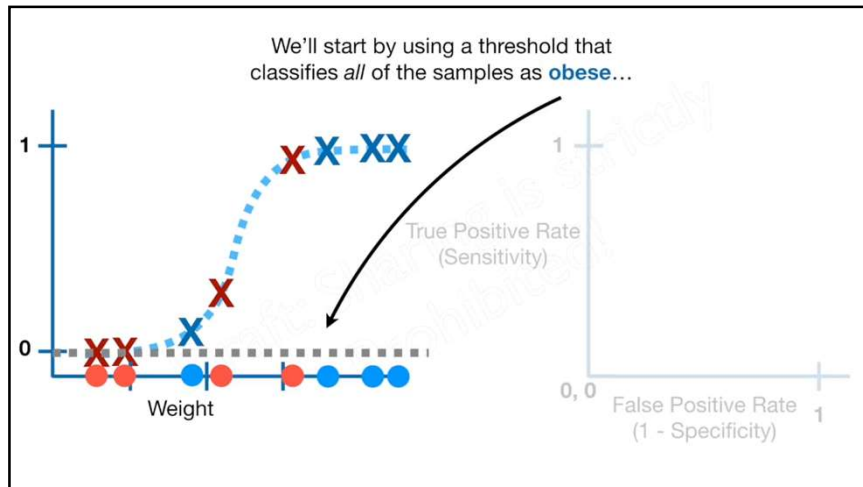


$$\text{True Positive Rate} = \text{Sensitivity} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

		Actual	
		Is Obese	Is Not Obese
Predicted	Is Obese	True Positives	False Positives
	Is Not Obese	False Negatives	True Negatives

The **True Positive Rate** tells you what proportion of **obese** samples were correctly classified.





$$\text{True Positive Rate} = \text{Sensitivity} = \frac{4}{4 + 0} = 1$$

$$\text{False Positive Rate} = 1 - \text{Specificity} = \frac{4}{4 + 0} = 1$$

The **False Positive Rate**, when the threshold is so low that every single sample is classified as **obese**, is also 1.

		Actual	
		Is Obese	Is Not Obese
Predicted	Is Obese	4	4
	Is Not Obese	0	0

$$\text{True Positive Rate} = \text{Sensitivity} = \frac{4}{4 + 0} = 1$$

$$\text{False Positive Rate} = 1 - \text{Specificity} = \frac{4}{4 + 0} = 1$$

This means that every single sample that was **not obese** was *incorrectly* classified as **obese**.

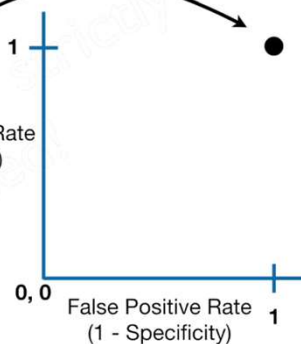
		Actual	
		Is Obese	Is Not Obese
Predicted	Is Obese	4	4
	Is Not Obese	0	0

$$\text{True Positive Rate} = \text{Sensitivity} = \frac{4}{4 + 0} = 1$$

$$\text{False Positive Rate} = 1 - \text{Specificity} = \frac{4}{4 + 0} = 1$$

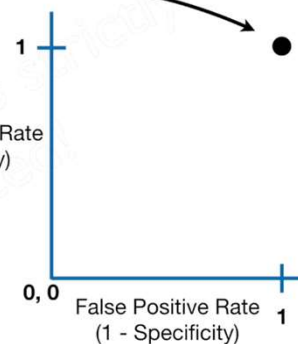
Now plot a point at 1, 1.

True Positive Rate (Sensitivity)

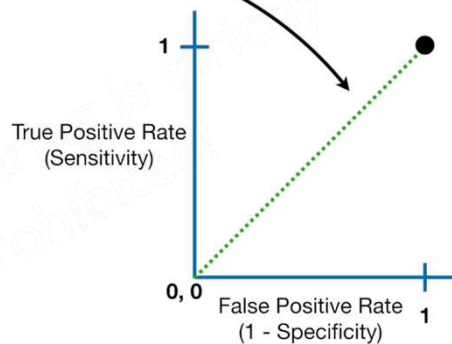


A point at 1, 1 means that even though we correctly classified **all** of the **obese** samples, we *incorrectly* classified **all** of the samples that were **not obese**.

True Positive Rate (Sensitivity)

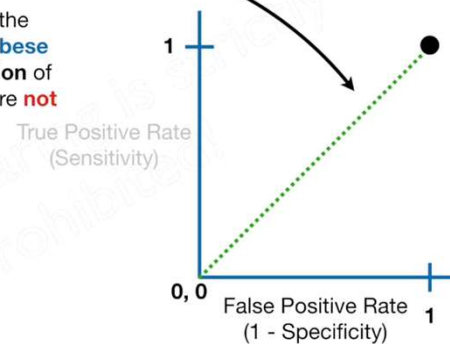


This **green diagonal line** shows where the **True Positive Rate = False Positive Rate**

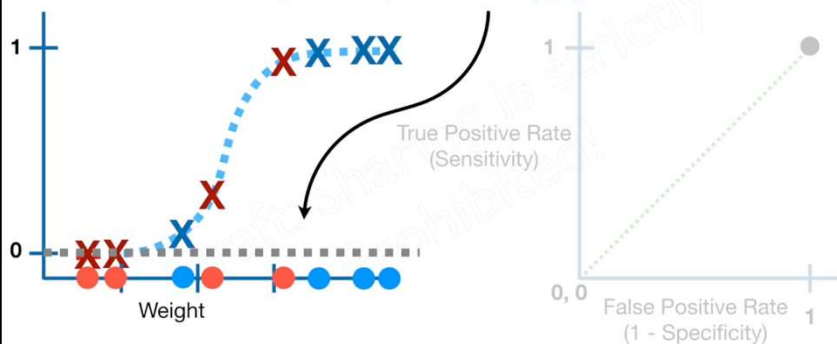


This **green diagonal line** shows where the **True Positive Rate = False Positive Rate**

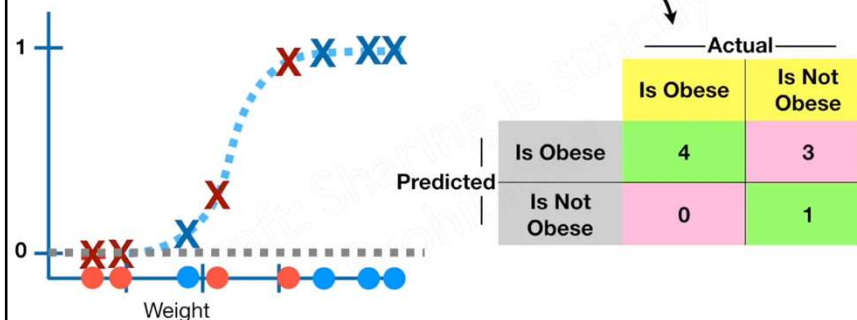
Any point on this **line** means that the **proportion of correctly classified obese** samples is the same as the **proportion of incorrectly classified samples that are not obese**.



Going back to the Logistic Regression, let's increase the threshold so that *all but the lightest sample* are called **obese**.



The new threshold gives us this **Confusion Matrix**.



True Positive Rate = Sensitivity = $\frac{4}{4+0} = 1$

False Positive Rate = 1 - Specificity = $\frac{3}{3+1} = 0.75$

True Positive Rate
(Sensitivity)

False Positive Rate
(1 - Specificity)

...and plot a
point at 0.75, 1.

Since the new point (0.75, 1) is to the left of the dotted green line, we know that the proportion of correctly classified samples that were obese (true positives) is greater than the proportion of the samples that were incorrectly classified as obese (false positives).

True Positive Rate
(Sensitivity)

False Positive Rate
(1 - Specificity)

In other words, the new threshold for deciding if a sample is obese or not is better than the first one.

True Positive Rate
(Sensitivity)

False Positive Rate
(1 - Specificity)

Now let's increase the threshold so that all but the two lightest samples are called obese.

obese.

True Positive Rate
(Sensitivity)

False Positive Rate
(1 - Specificity)

Weight

True Positive Rate = Sensitivity = $\frac{4}{4+0} = 1$

False Positive Rate = 1 - Specificity = $\frac{2}{2+2} = 0.5$

...and plot a point at 0.5, 1.

True Positive Rate
(Sensitivity)

False Positive Rate
(1 - Specificity)

The new point (0.5, 1) is even further to the left of the dotted green line, showing that the new threshold further decreases the proportion of the samples that were incorrectly classified as obese (false positives).

True Positive Rate
(Sensitivity)

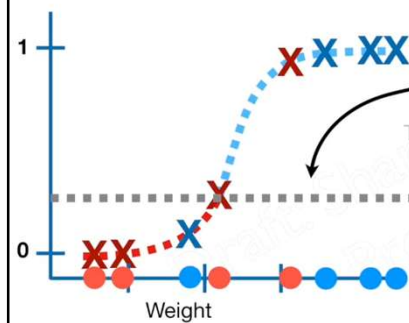
False Positive Rate
(1 - Specificity)

In other words, the new threshold is the best one so far.

True Positive Rate
(Sensitivity)

False Positive Rate
(1 - Specificity)

Now we increase the threshold again...



True Positive Rate
(Sensitivity)

False Positive Rate
(1 - Specificity)

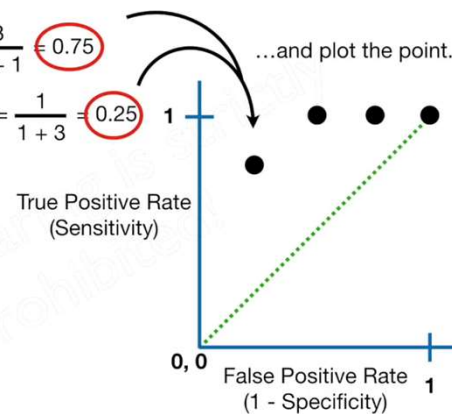
True Positive Rate = Sensitivity = $\frac{3}{3+1} = 0.75$...calculate the **True Positive Rate** and the **False Positive Rate**...

False Positive Rate = 1 - Specificity = $\frac{1}{1+3} = 0.25$

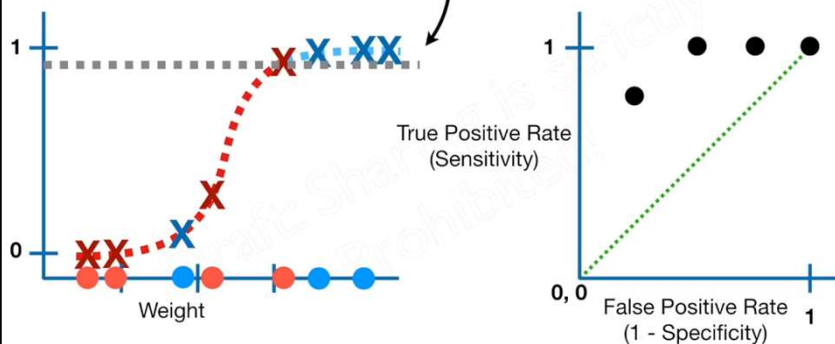
		Actual	
		Is Obese	Is Not Obese
Predicted	Is Obese	3	1
	Is Not Obese	1	3

True Positive Rate = Sensitivity = $\frac{3}{3+1} = 0.75$

False Positive Rate = 1 - Specificity = $\frac{1}{1+3} = 0.25$



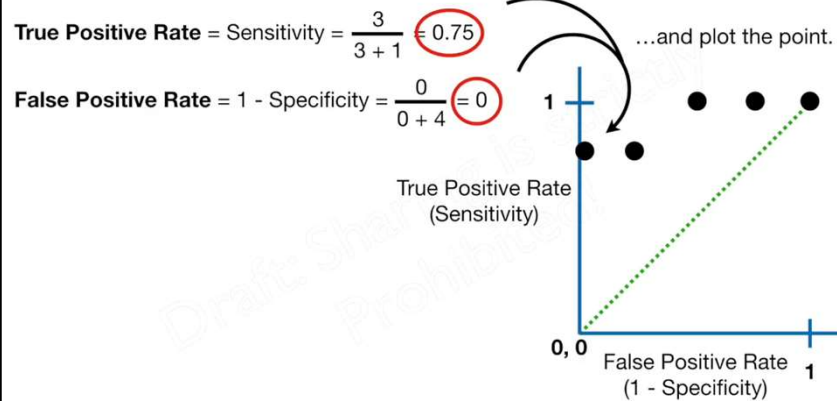
Now we increase the threshold again...



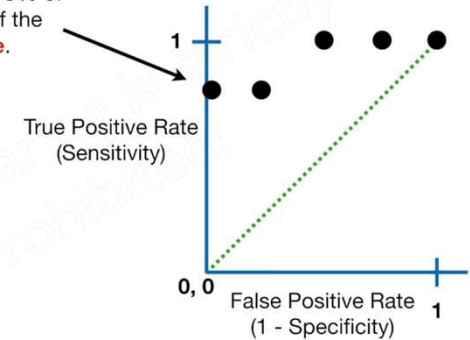
True Positive Rate = Sensitivity = $\frac{3}{3+1} = 0.75$...calculate the **True Positive Rate** and the **False Positive Rate**...

False Positive Rate = 1 - Specificity = $\frac{0}{0+4} = 0$

		Actual	
		Is Obese	Is Not Obese
Predicted	Is Obese	3	0
	Is Not Obese	1	4

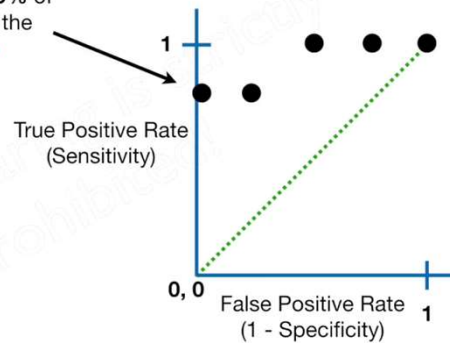


The threshold represented by the new point (0, 0.75) correctly classified 75% of the **obese** samples and 100% of the samples that were **not obese**.

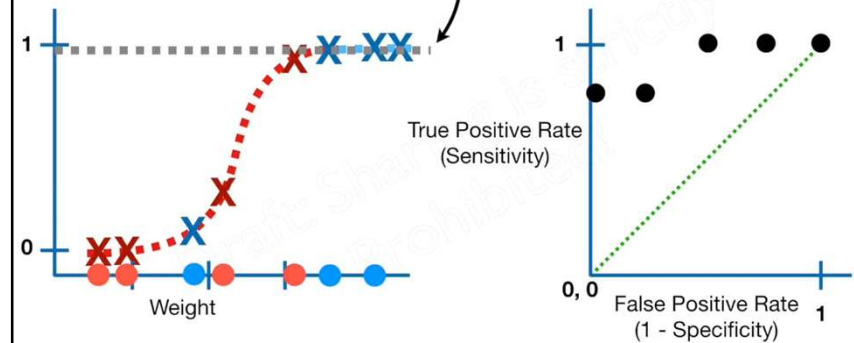


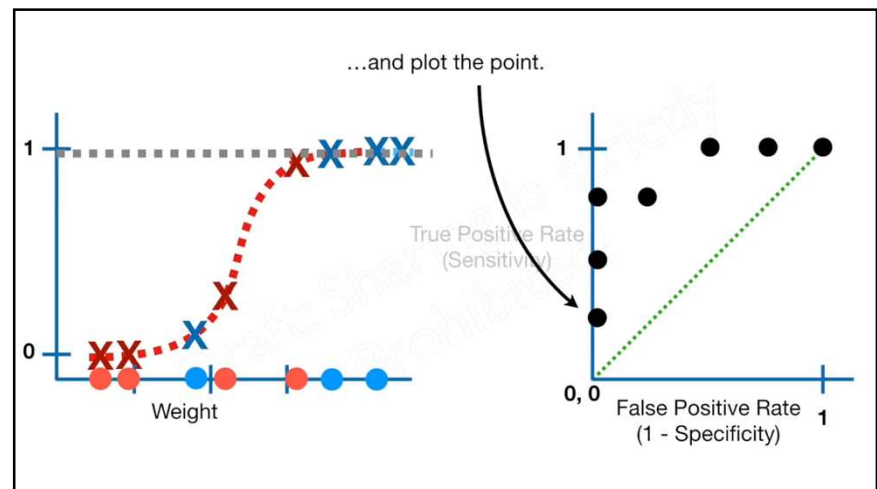
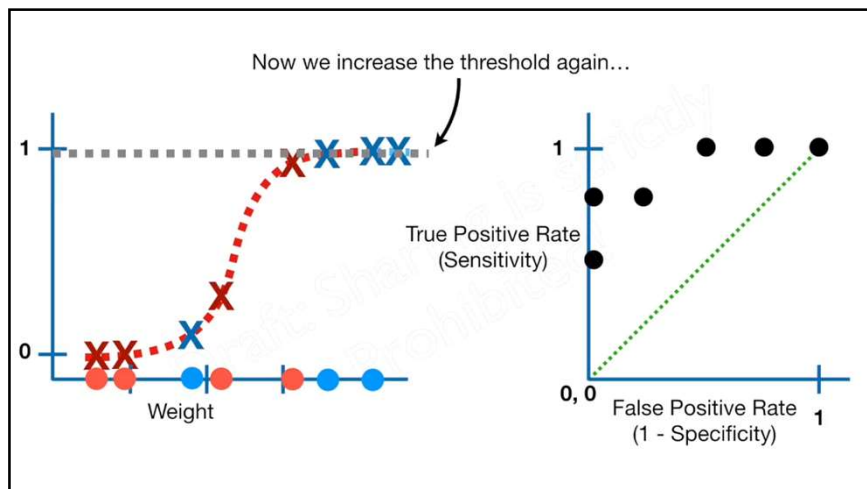
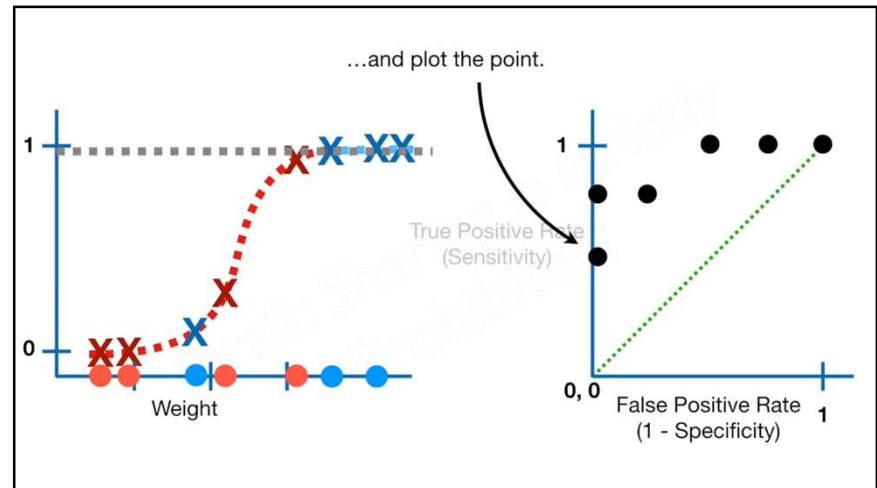
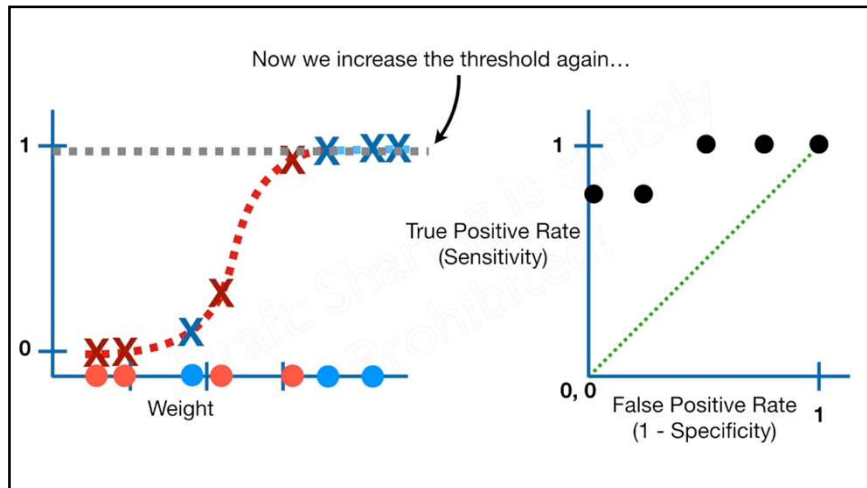
The threshold represented by the new point (0, 0.75) correctly classified 75% of the **obese** samples and 100% of the samples that were **not obese**.

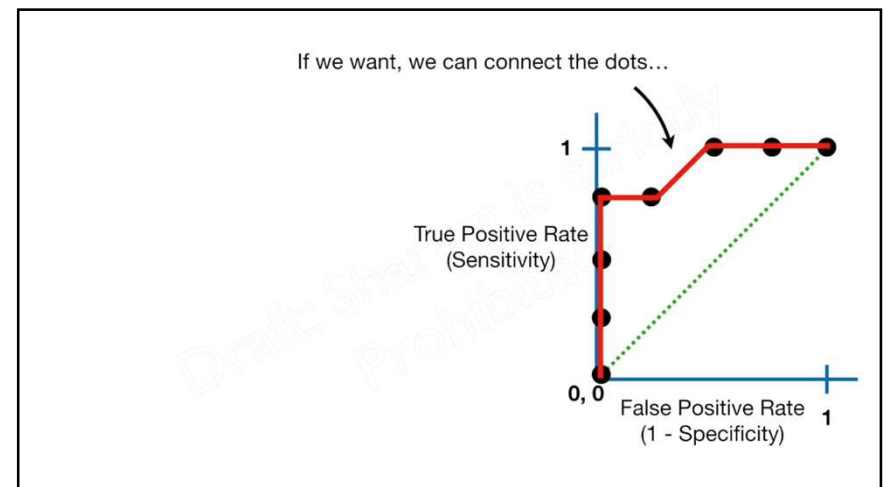
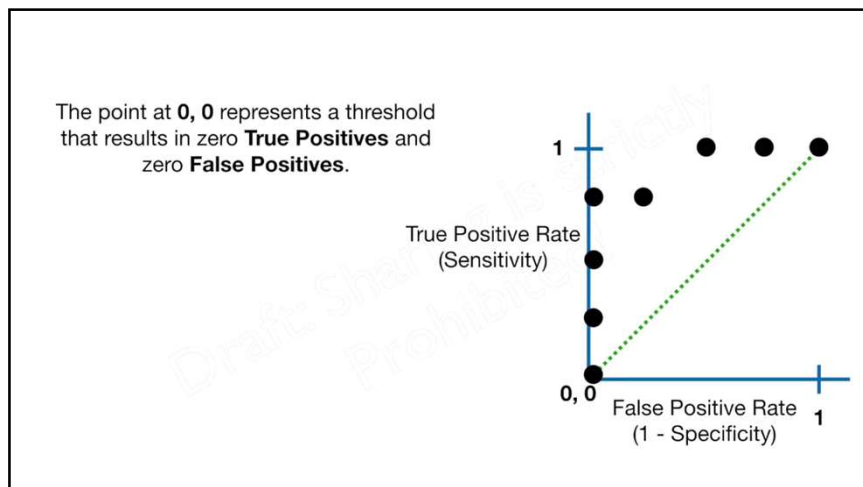
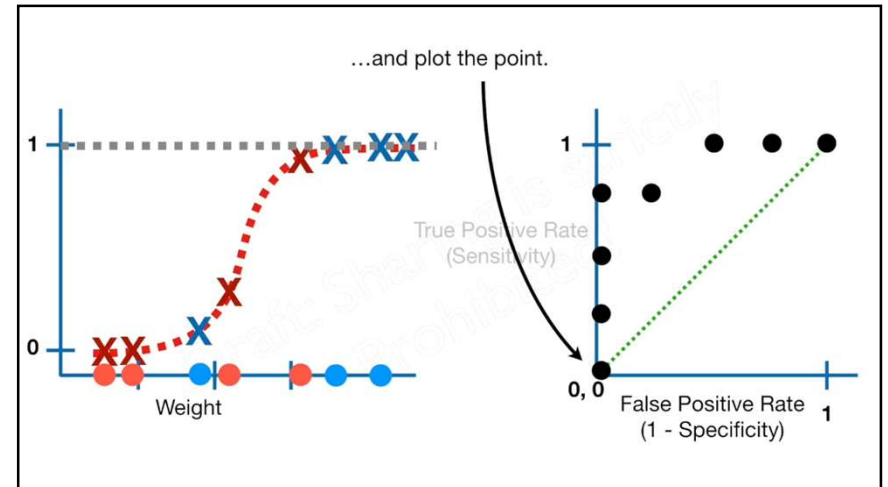
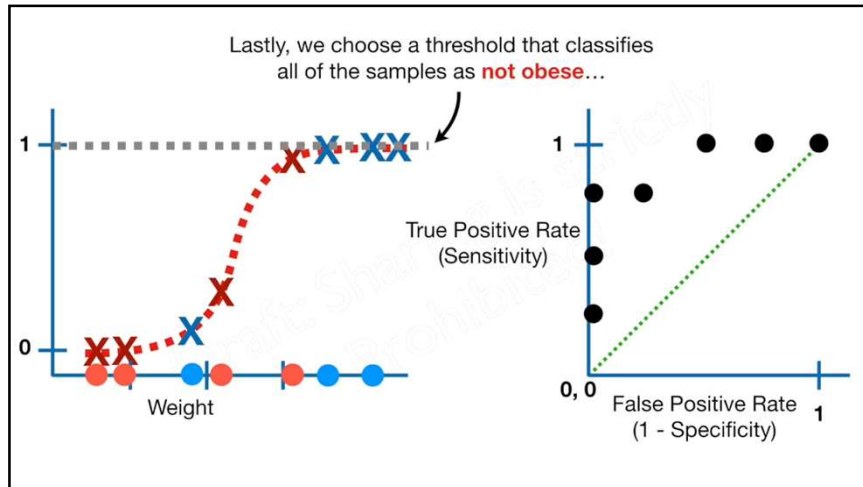
In other words, this threshold resulted in no **False Positives**.



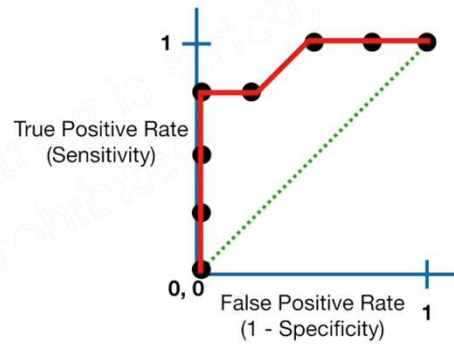
Now we increase the threshold again...



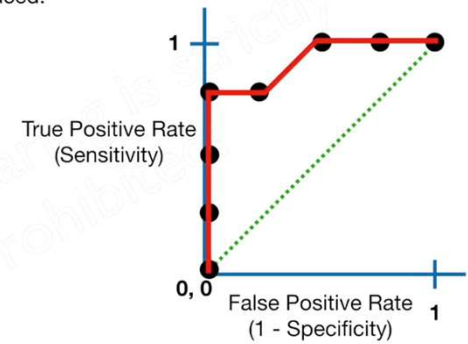




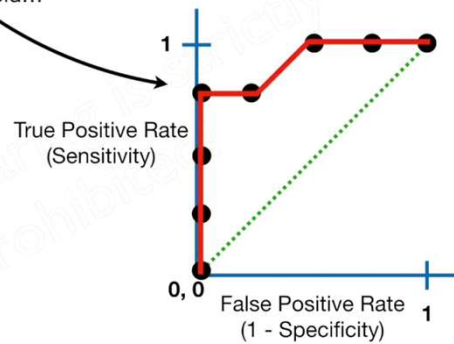
...and that gives us an **ROC** graph.



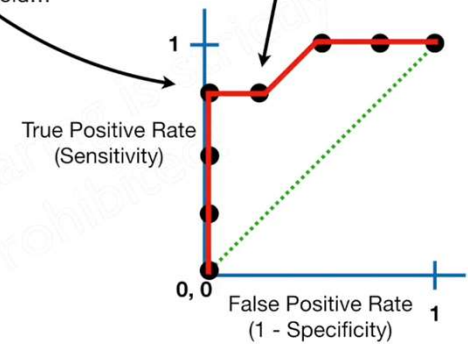
The **ROC** graph summarizes all of the confusion matrices that each threshold produced.



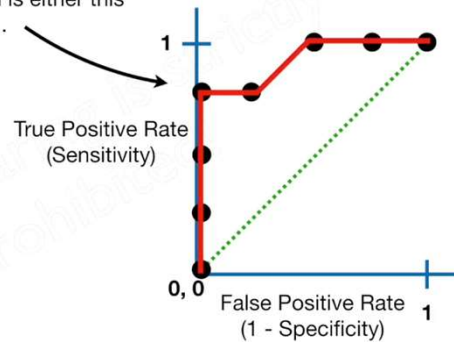
Without having to sort through the confusion matrices, I can tell that this threshold...



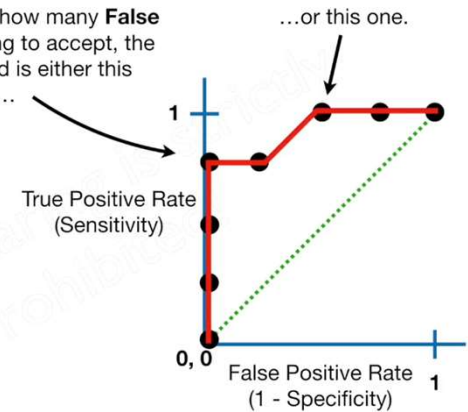
Without having to sort through the confusion matrices, I can tell that this threshold...
...is better than this threshold.



And depending on how many **False Positives** I'm willing to accept, the optimal threshold is either this one...

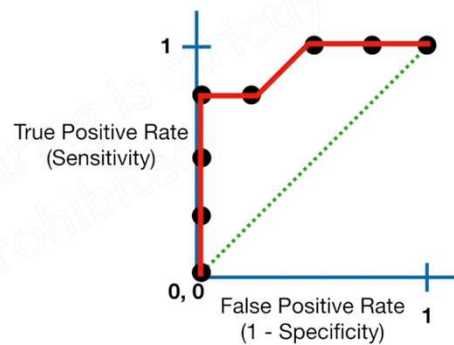


And depending on how many **False Positives** I'm willing to accept, the optimal threshold is either this one...

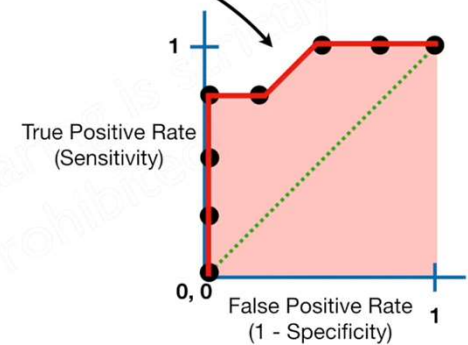


...or this one.

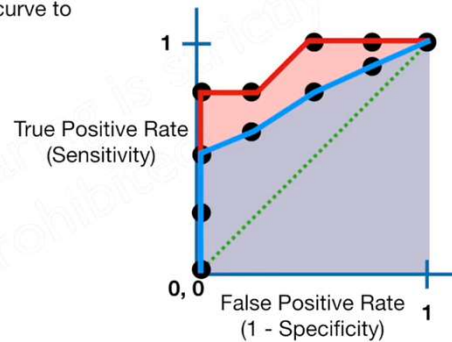
Now that we know what an **ROC** graph is, let's talk about the **Area Under the Curve**, or **AUC**...



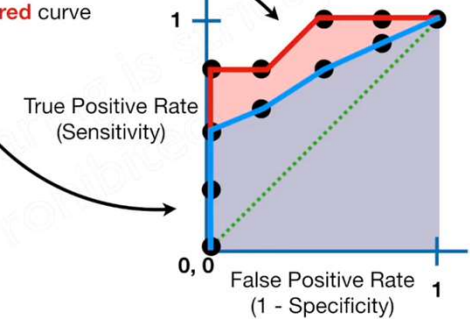
The **AUC** (Area Under the Curve) is **0.9**



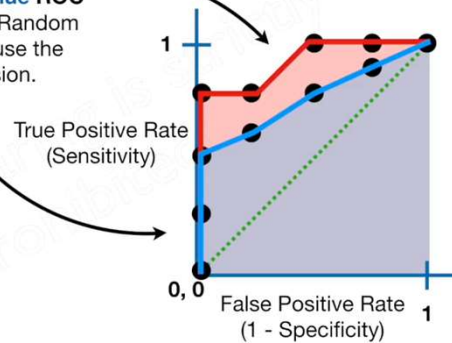
The **AUC** makes it easy to compare one **ROC** curve to another.



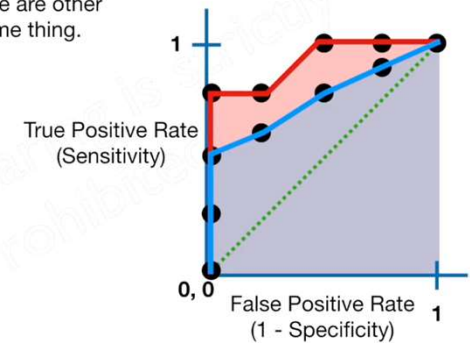
The **AUC** for the **red ROC** curve is greater than the **AUC** for the **blue ROC** curve, suggesting that the **red** curve is better.



So if the **red ROC** curve represented Logistic Regression and the **blue ROC** curve represented a Random Forest, you would use the Logistic Regression.

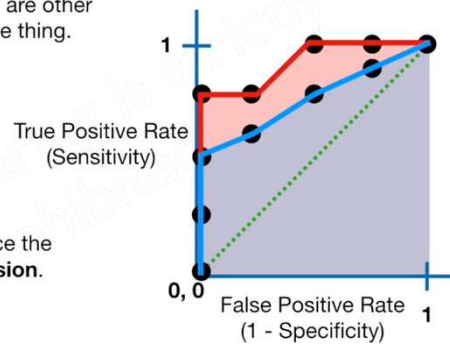


Although **ROC** graphs are drawn using **True Positive Rates** and **False Positive Rates** to summarize confusion matrices, there are other metrics that attempt to do the same thing.



Although **ROC** graphs are drawn using **True Positive Rates** and **False Positive Rates** to summarize confusion matrices, there are other metrics that attempt to do the same thing.

For example, people often replace the **False Positive Rate** with **Precision**.

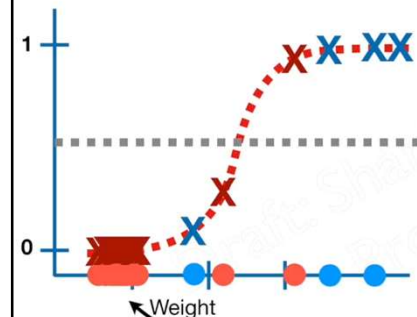


$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

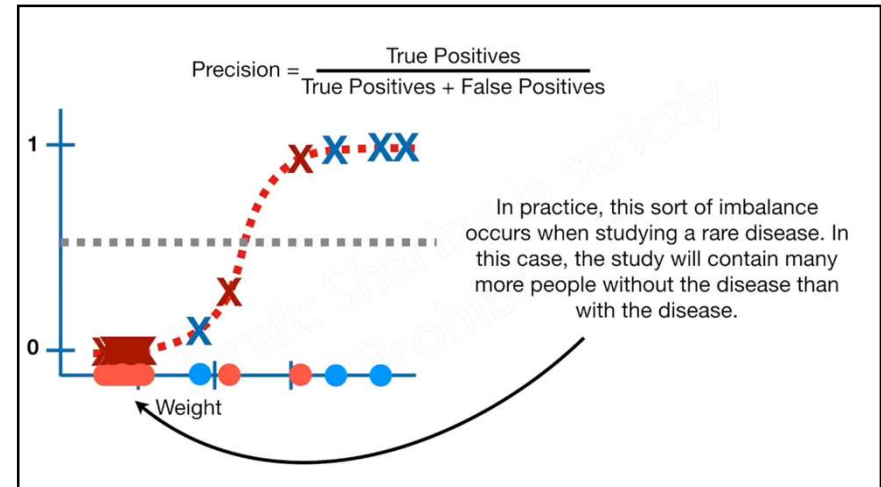
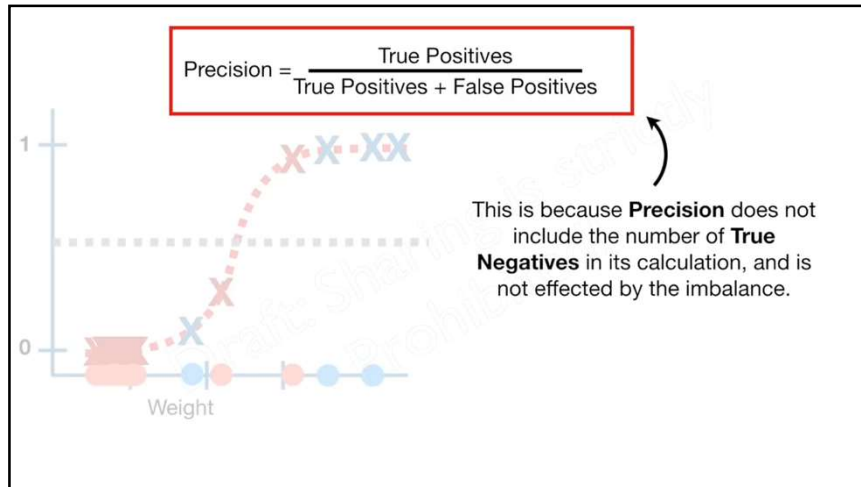
		Actual	
		Is Obese	Is Not Obese
Predicted	Is Obese	True Positives	False Positives
	Is Not Obese	False Negatives	True Negatives

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}} = \text{Precision is the proportion of positive results that were correctly classified}$$

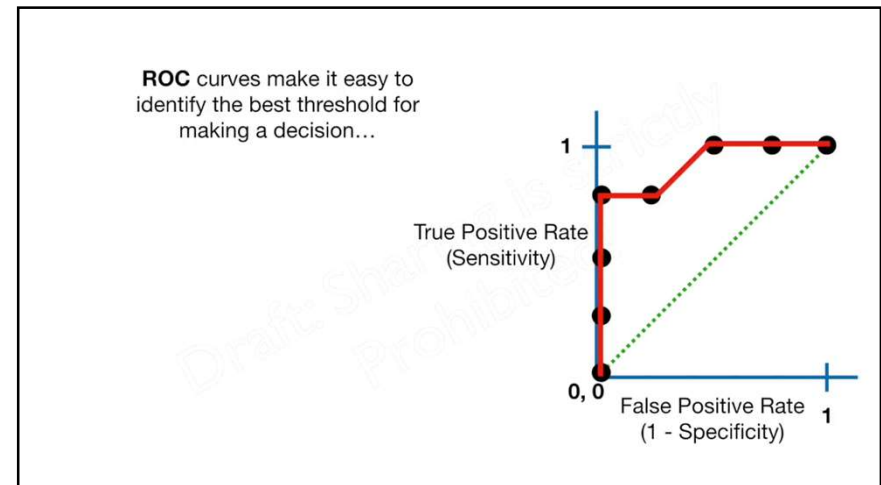
		Actual	
		Is Obese	Is Not Obese
Predicted	Is Obese	True Positives	False Positives
	Is Not Obese	False Negatives	True Negatives

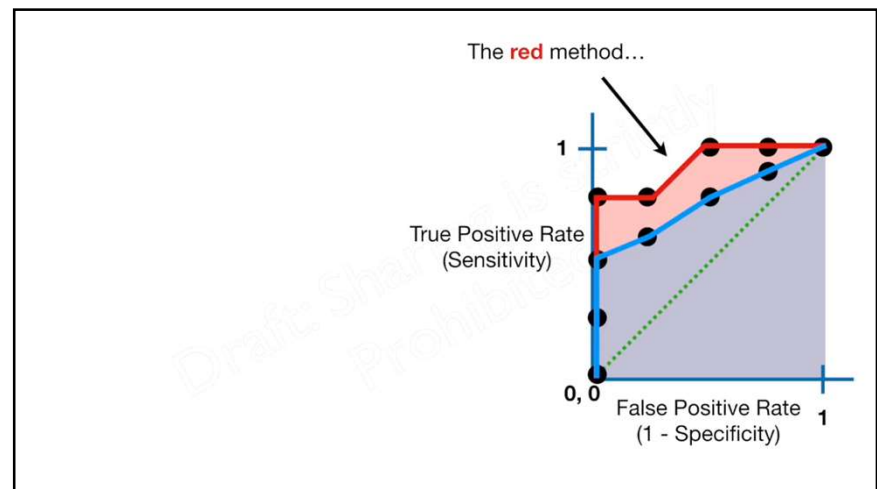
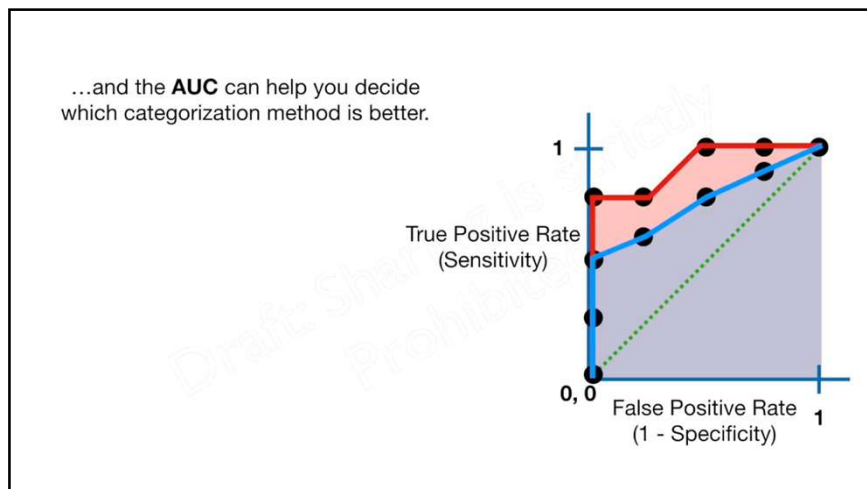
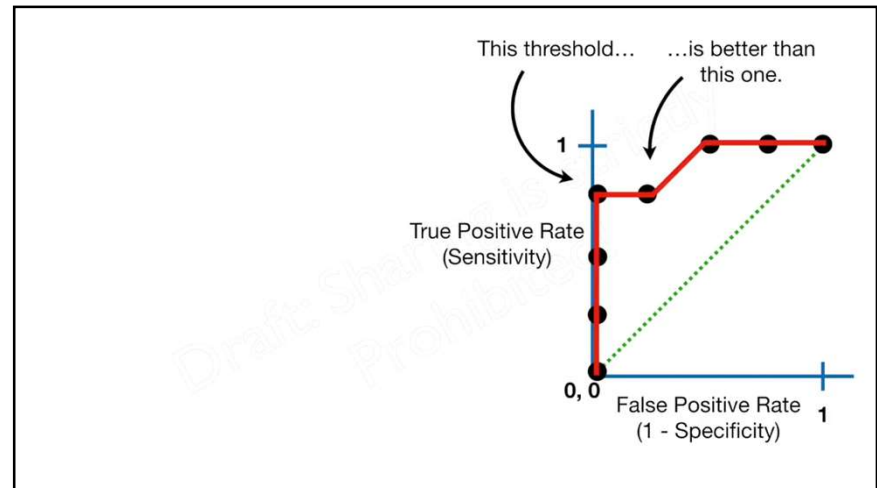
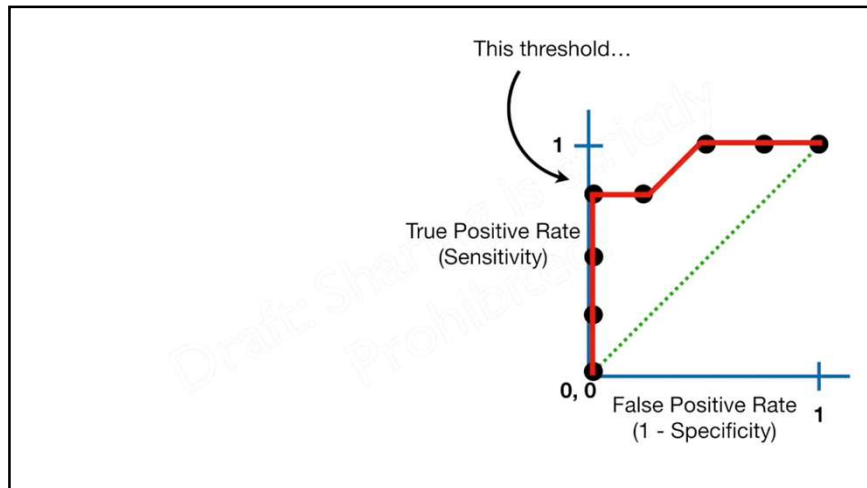


If there were lots of samples that were **not obese** relative to the number of **obese** samples, then **Precision** might be more useful than the **False Positive Rate**.



In Summary...





THANK YOU!