

NYU FRE 7773 - Week 4

Machine Learning in Financial Engineering
Ethan Rosenthal

Linear Models for Classification

Machine Learning in Financial Engineering
Ethan Rosenthal

Classification

Announcing \$0 Online Stock, ETF, and Options Commissions

[GET THE DETAILS](#)

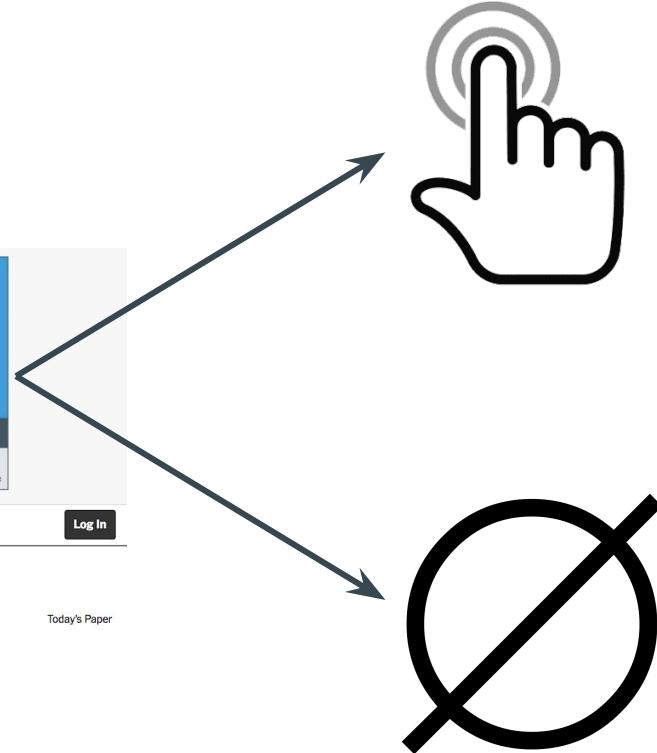
+ A Satisfaction Guarantee
If you're not happy for any reason, we'll refund your eligible fee and work with you to make things right.

charles schwab
Own your tomorrow

≡ Q ENGLISH ESPAÑOL 中文 Log In

The New York Times

Today's Paper





Share



Share

<https://www.theverge.com/2017/6/26/15876006/hot-dog-app-android-silicon-valley>

<https://huggingface.co/spaces/huggingface-projects/llama-2-7b-chat>

Logistic Regression

The ML Recipe

1. Think up some model
2. Feed **data** into the model and make predictions.
3. Calculate the loss between predictions and true values.
4. Determine the model parameters that produce the minimum loss.

The ML Recipe

1. Think up some model
2. Feed **data** into the model and make predictions.
3. Calculate the loss between predictions and true values.
4. Determine the model parameters that produce the minimum loss.

Logistic Regression – The Model

- Recall our linear model for regression

$$\hat{y}_i = \sum_j^p x_{ij} \cdot \beta_j$$

- We can't use this for classification because we must predict 0 or 1.
- So, let's squash this between 0 and 1.

Logistic Regression – The Model

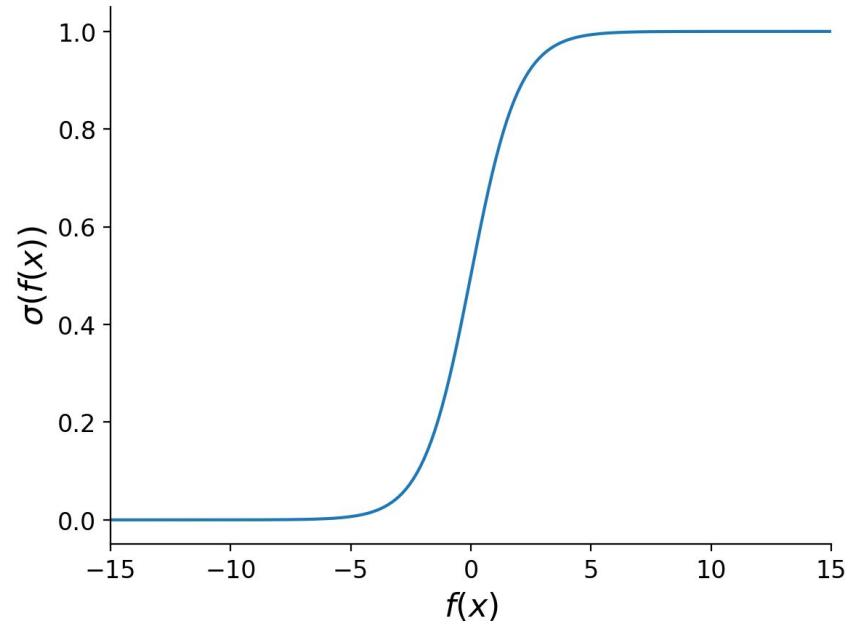
- Squash the linear model between 0 and 1.

$$\hat{y}_i = \sum_j^p x_{ij} \cdot \beta_j$$

$$f(x_i) = \sum_j^p x_{ij} \cdot \beta_j$$

$$\hat{y}_i = \sigma(f(x_i))$$

$$\sigma(f(x_i)) = \frac{1}{1 + e^{-f(x_i)}}$$

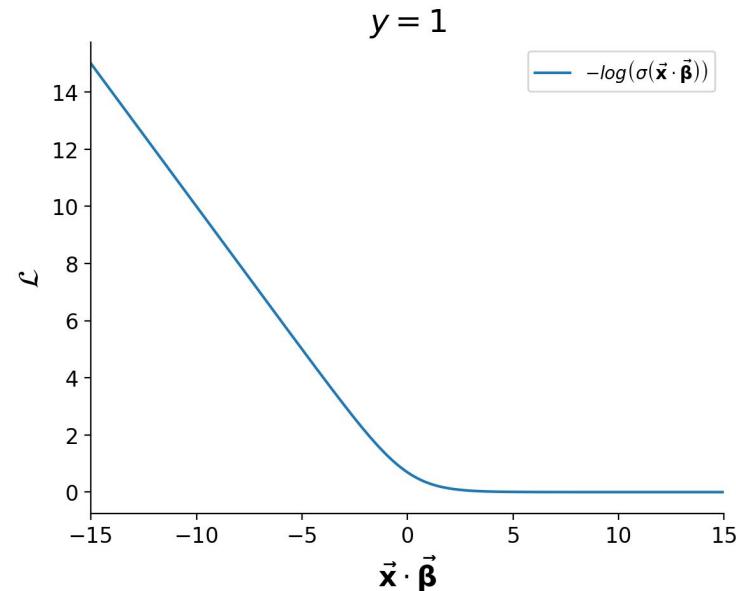


The ML Recipe

1. Think up some model
2. Feed **data** into the model and make predictions.
3. Calculate the loss between predictions and true values.
4. Determine the model parameters that produce the minimum loss.

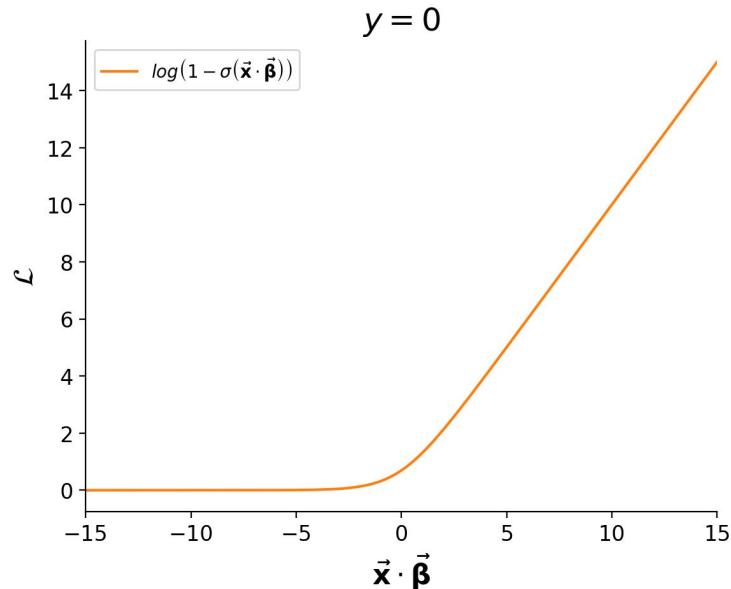
Logistic Regression – The Loss

$$\mathcal{L}(\vec{\beta}) = - \sum_{i=1}^N y_i \log \left(\sigma(\vec{x}_i \cdot \vec{\beta}) \right)$$



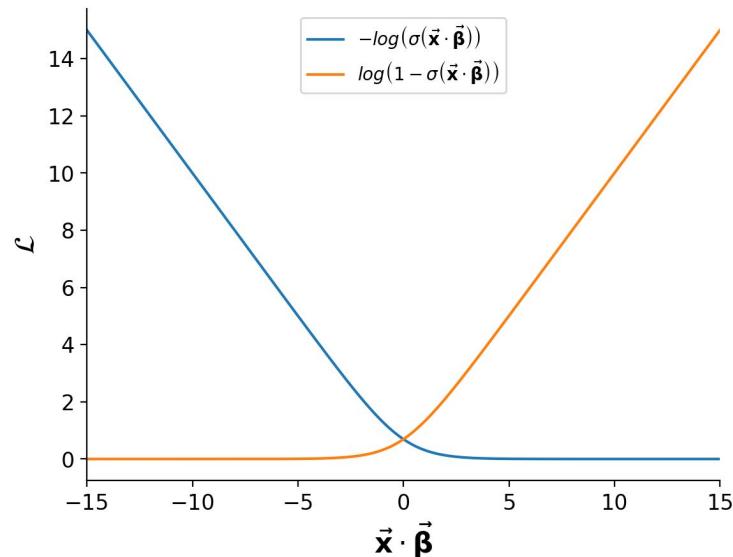
Logistic Regression – The Loss

$$\mathcal{L}(\vec{\beta}) = - \sum_{i=1}^N + (1 - y_i) \log \left(1 - \sigma(\vec{x}_i \cdot \vec{\beta}) \right)$$



Logistic Regression – The Loss

$$\mathcal{L}(\vec{\beta}) = - \sum_{i=1}^N y_i \log \left(\sigma(\vec{x}_i \cdot \vec{\beta}) \right) + (1 - y_i) \log \left(1 - \sigma(\vec{x}_i \cdot \vec{\beta}) \right)$$



The ML Recipe

1. Think up some model
2. Feed **data** into the model and make predictions.
3. Calculate the loss between predictions and true values.
4. Determine the model parameters that produce the minimum loss.

The ML Recipe

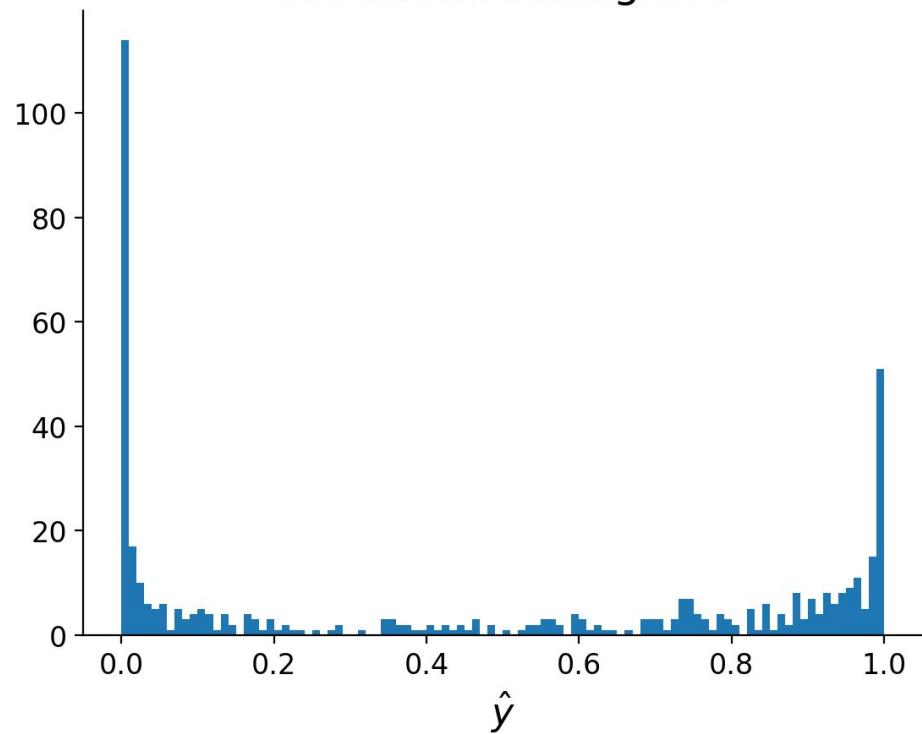
1. Think up some model
2. Feed **data** into the model and make predictions.
3. Calculate the loss between predictions and true values.
4. Determine the model parameters that produce the minimum loss.

$$\frac{\partial \mathcal{L}}{\partial \vec{\beta}} = \sum_{i=1}^N \left(y_i - \frac{1}{1 + e^{(\vec{x}_i \cdot \vec{\beta})}} \right) \vec{x}_i$$

Logistic Regression – The Predictions

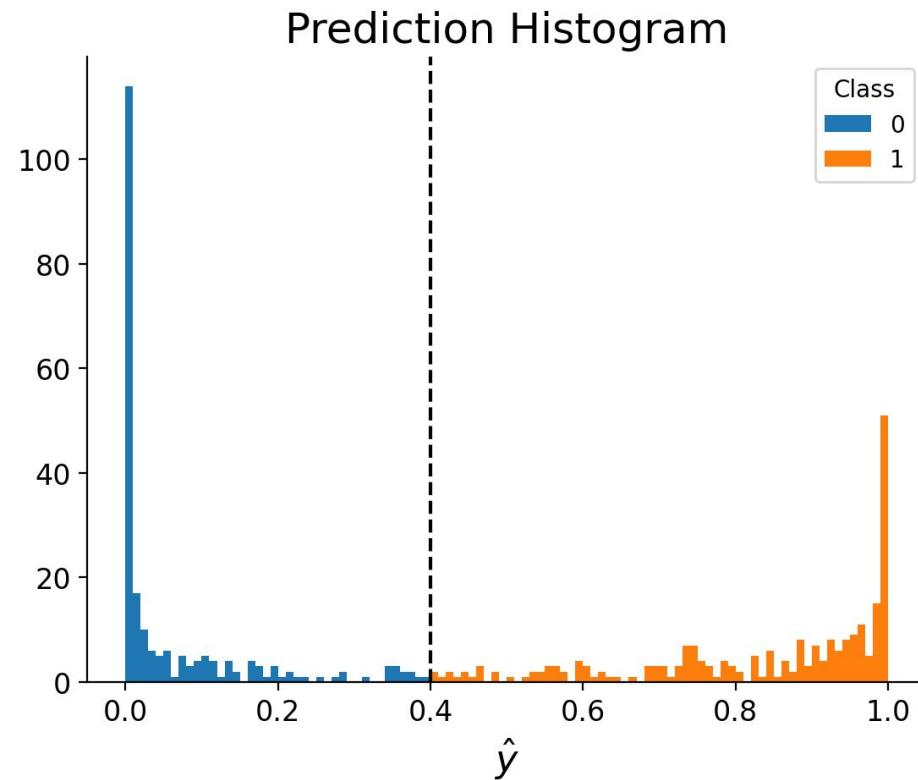
- Predictions lie between 0 and 1.
- They can be interpreted as probabilities.

Prediction Histogram



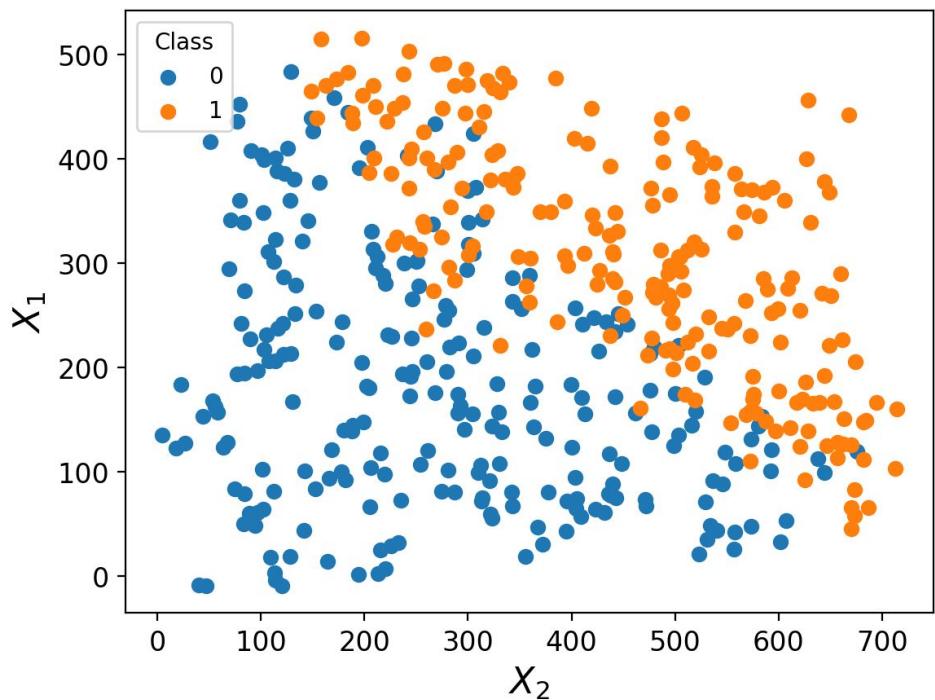
Logistic Regression – The Predictions

- Predictions lie between 0 and 1.
- They can be interpreted as probabilities.
- We pick a **threshold** to classify probabilities into classes.



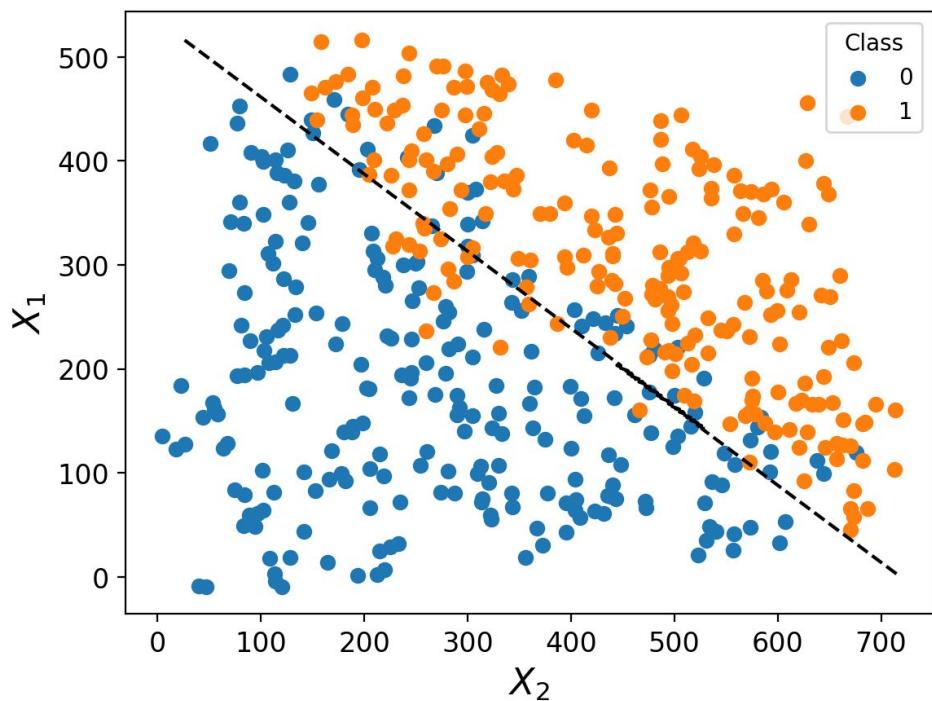
Logistic Regression – The Decision Boundary

- Imagine a classification problem with 2 features.
- For any value of X_1 and X_2 , our model predicts a probability between 0 and 1.
- Predictions above (below) the threshold are predicted as 1 (0).



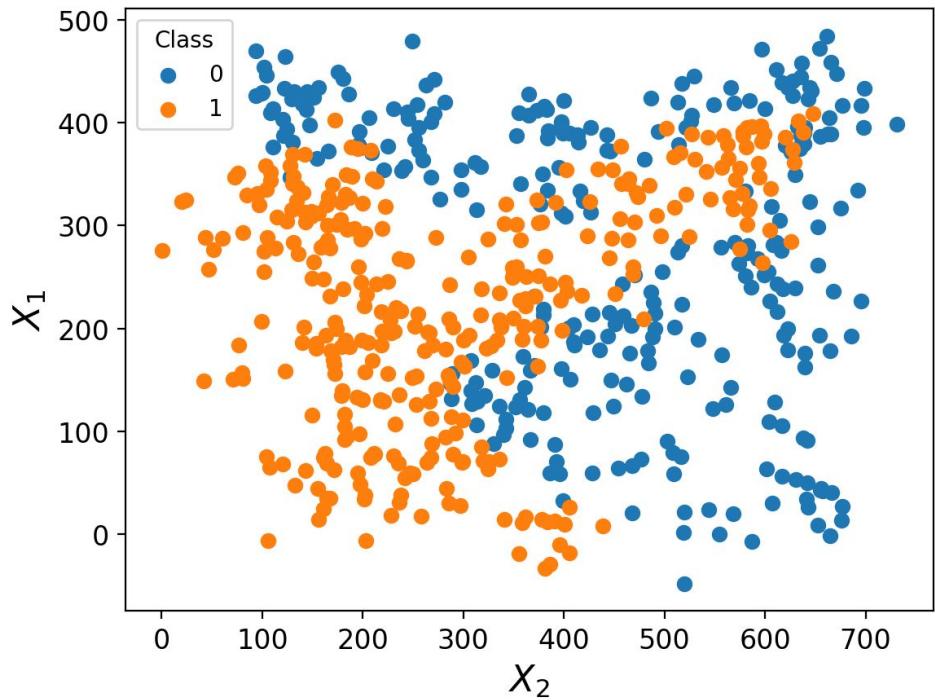
Logistic Regression – The Decision Boundary

- Imagine a classification problem with 2 features.
- For any value of X_1 and X_2 , our model predicts a probability between 0 and 1.
- Predictions above (below) the threshold are predicted as 1 (0).
- The line that runs along the region of predictions where the prediction == threshold is the **decision boundary**.



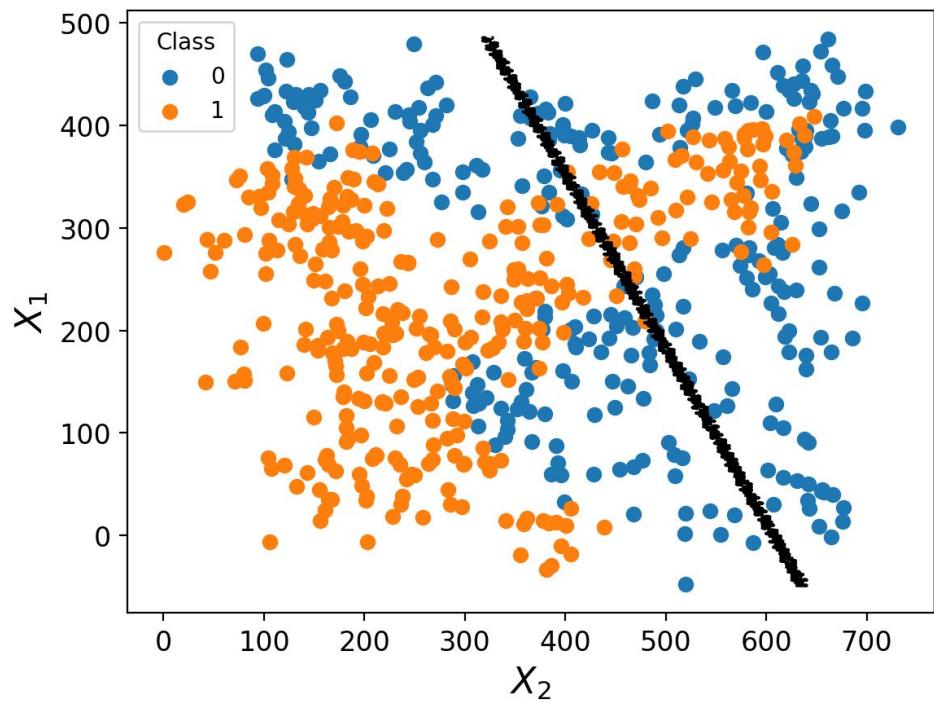
Logistic Regression – The Decision Boundary

- Problems are often nonlinear.



Logistic Regression – The Decision Boundary

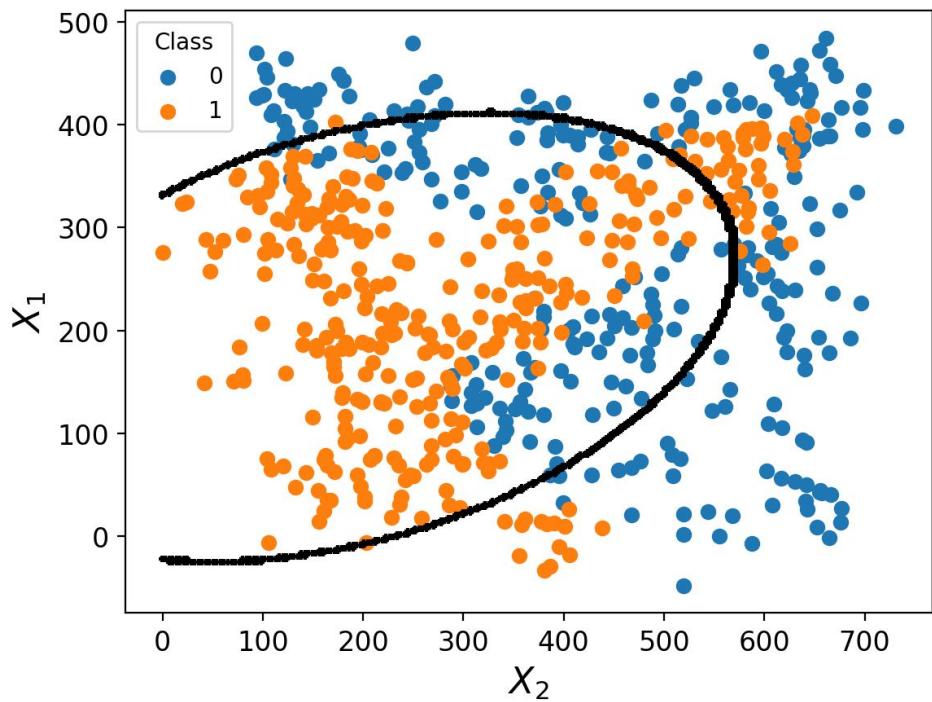
- Problems are often nonlinear.
- And linear decision boundaries don't work well on them!



Logistic Regression – The Decision Boundary

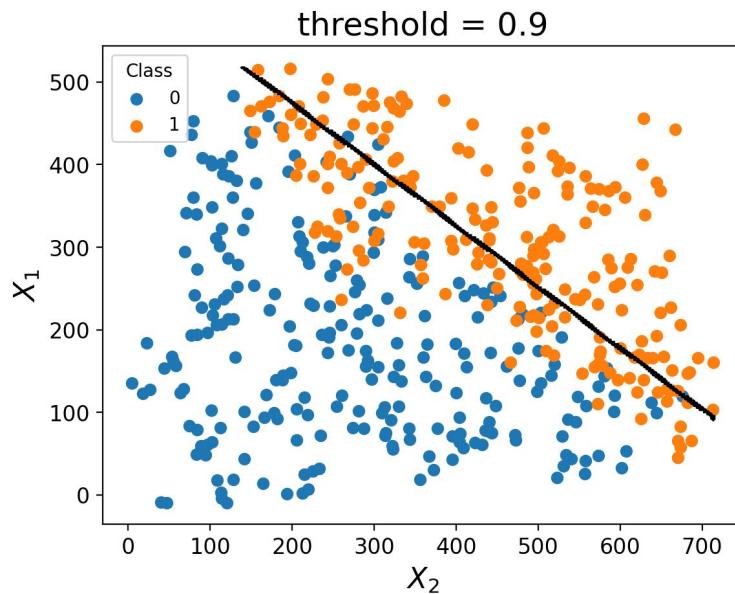
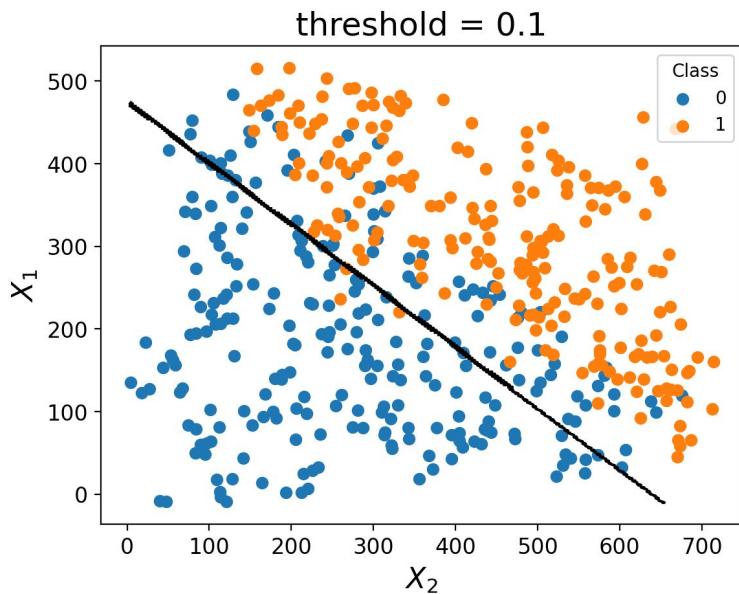
- Problems are often nonlinear.
- And linear decision boundaries don't work well on them!
- Can try creating nonlinear features:

$$\begin{aligned}\hat{y}_i &= \beta_0 + \beta_1 \cdot X_{i1} + \beta_2 \cdot X_{i2} \\ &+ \beta_3 \cdot X_{i1} \cdot X_{i2} + \beta_4 \cdot X_{i1}^2 \\ &+ \beta_5 \cdot X_{i2}^2\end{aligned}$$



Logistic Regression – The Decision Boundary

Your threshold / decision boundary is a choice!

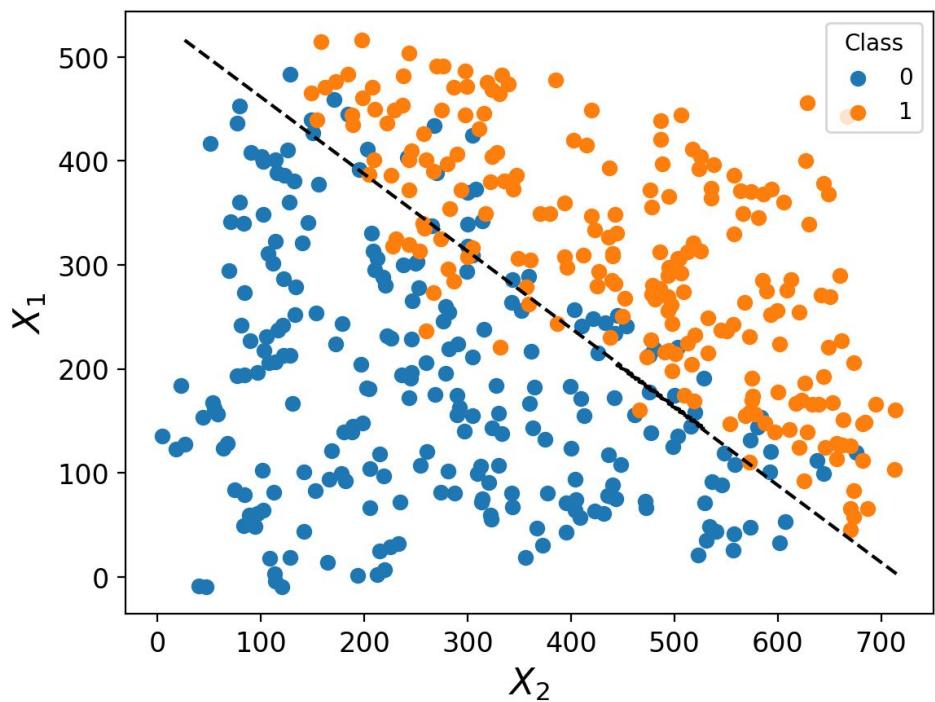


Classification Evaluation

Classification Evaluation

Confusion Matrix

		Predicted Labels
True Labels	False	True
	False	True
False	204	41
True	16	198

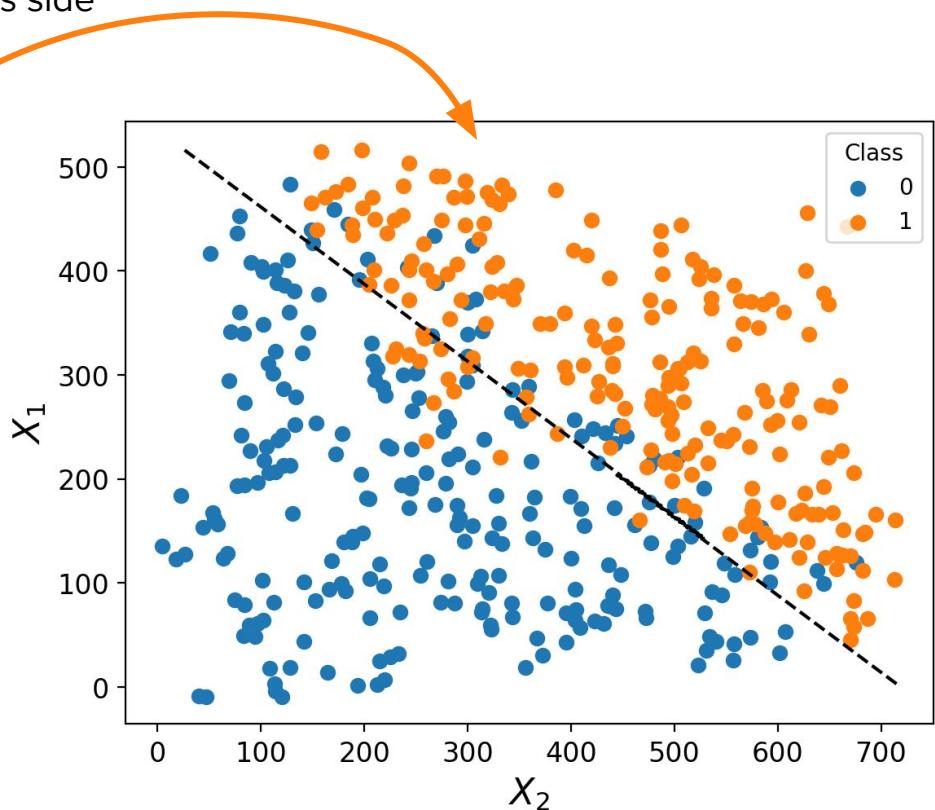


Classification Evaluation

Confusion Matrix

		Predicted Labels
True Labels	False	True
	False	True
False	204	41
True	16	198

True Positives
All orange dots on
this side

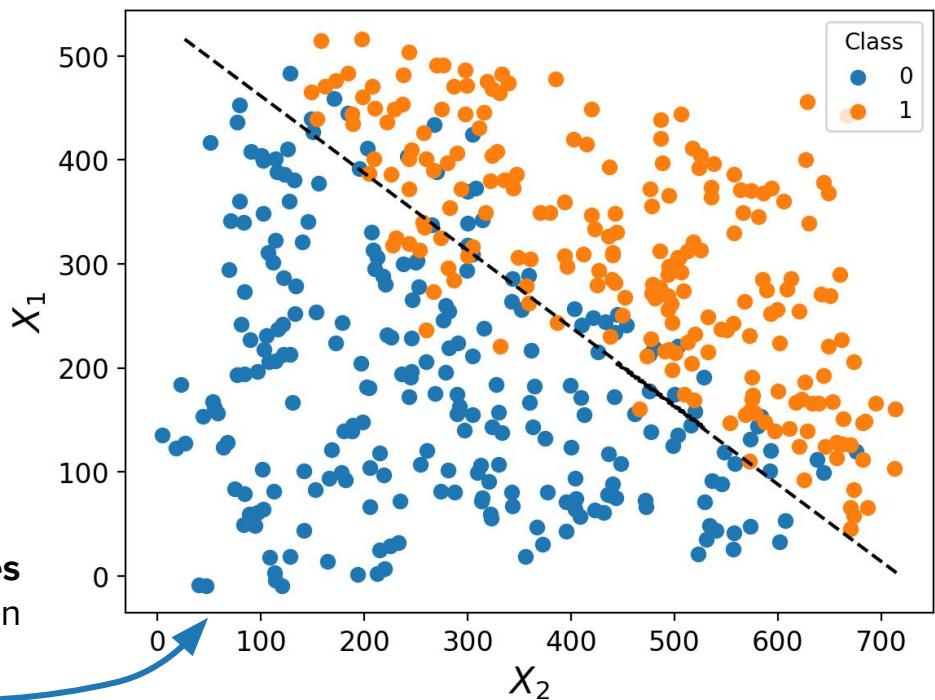


Classification Evaluation

Confusion Matrix

		Predicted Labels
True Labels	False	True
	False	True
False	204	41
True	16	198

True Negatives
All blue dots on
this side



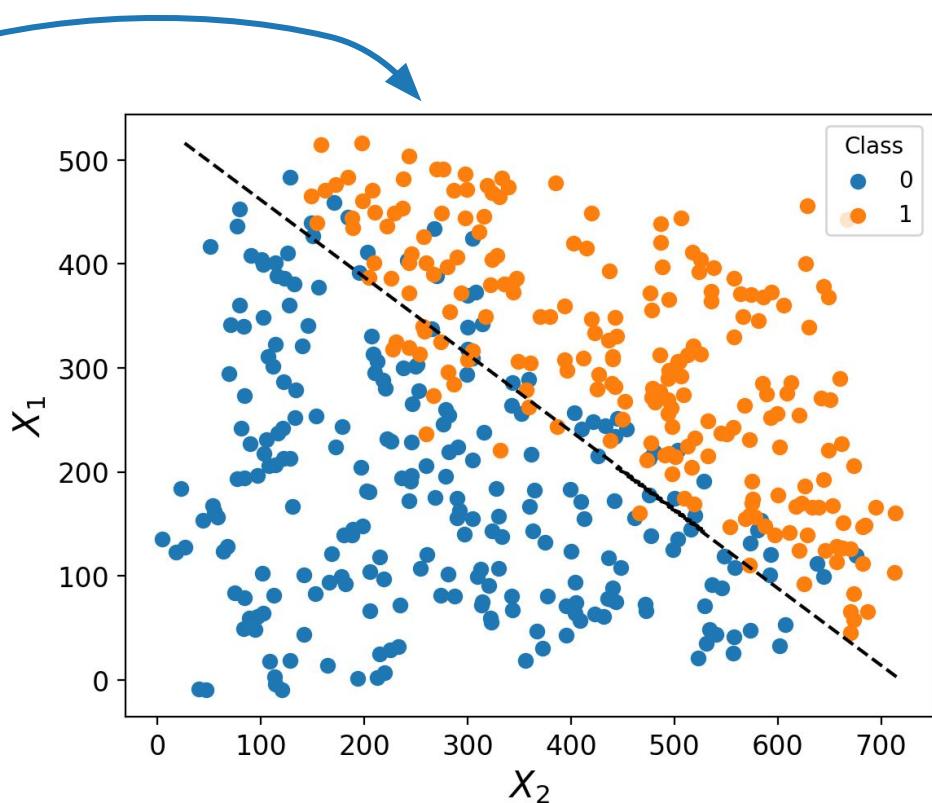
Classification Evaluation

Confusion Matrix

		Predicted Labels
True Labels	False	True
	False	True
False	204	41
True	16	198

False Positives

All blue dots on this side

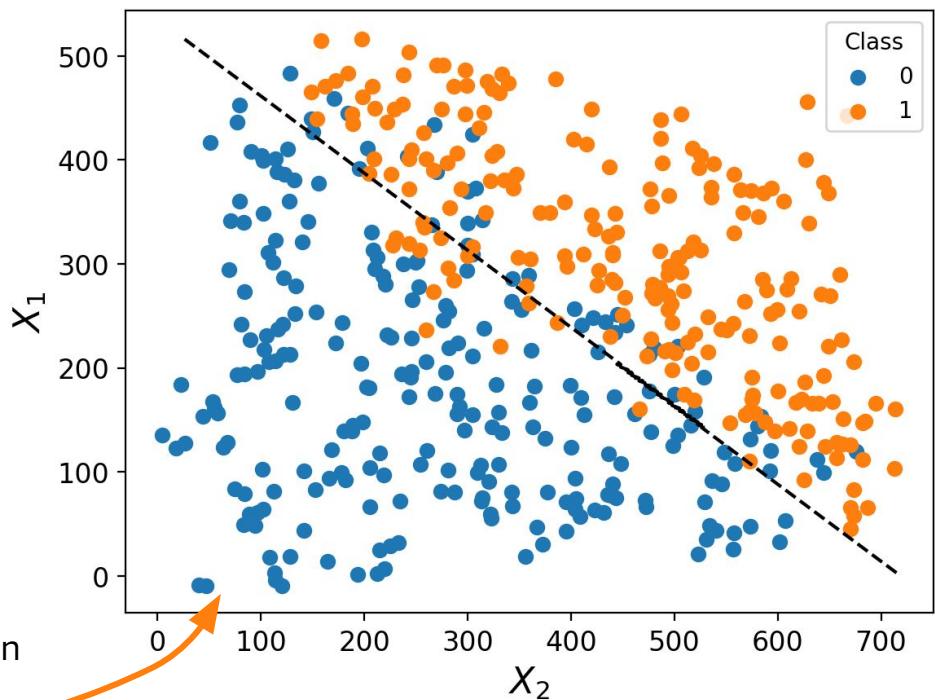


Classification Evaluation

Confusion Matrix

		Predicted Labels
True Labels	False	True
	False	True
False	204	41
True	16	198

False Negatives
All orange dots on
this side



Classification Evaluation

Confusion Matrix

		Predicted Labels
True Labels	False	True
	False	True
False	204	41
True	16	198

$$\text{accuracy} = \frac{TP + TN}{N}$$

Total samples →

$$\text{precision} = \frac{TP}{TP + FP}$$

$$\text{recall} = \frac{TP}{TP + FN}$$

Classification Evaluation

Confusion Matrix

		Predicted Labels
		False
True Labels	False	True
	204	41
True	16	198

$$precision = \frac{TP}{TP + FP}$$

$$recall = \frac{TP}{TP + FN}$$

Precision: For all of the samples my model predicted positive, what % of them were actually positive?

Recall: Out of all positive examples in my data, what % did my model predict positive?

Classification Evaluation

Confusion Matrix

		Predicted Labels	False	True
		True Labels	False	True
False	False	204	41	
	True	16		198

$$precision = \frac{TP}{TP + FP}$$

$$recall = \frac{TP}{TP + FN}$$

F1 Score: Harmonic mean of precision and recall

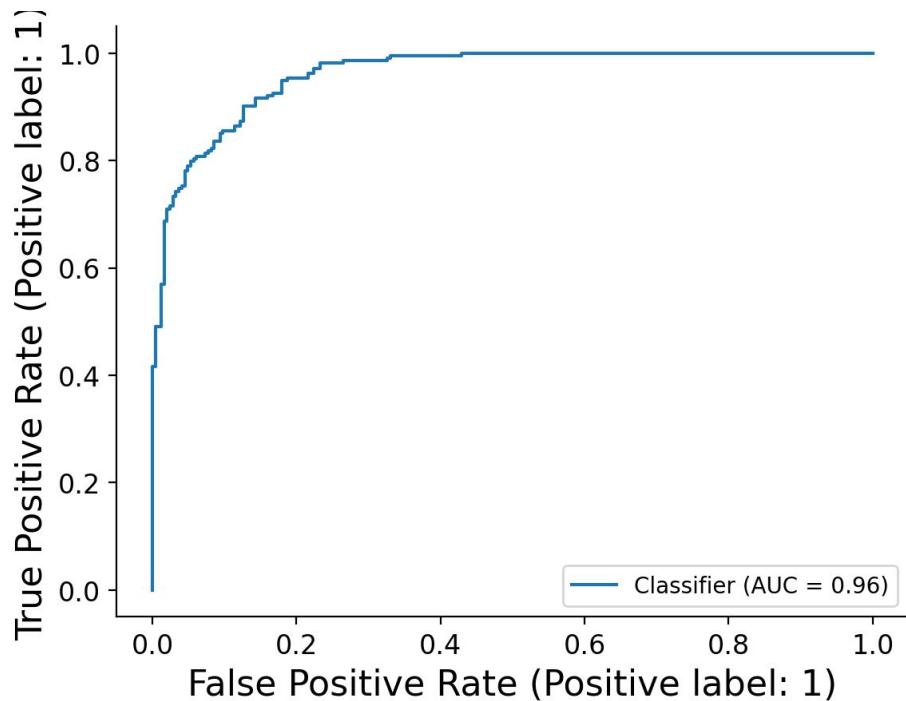
$$\begin{aligned} F_1 &= \frac{2}{\frac{1}{precision} + \frac{1}{recall}} \\ &= 2 \frac{precision \cdot recall}{precision + recall} \end{aligned}$$

The ML Recipe

1. Think up some model
2. Feed **data** into the model and make predictions.
3. Calculate the loss between predictions and true values.
4. Determine the model parameters that produce the minimum loss.
5. Pick a threshold if it's a classification model.
 - a. The threshold is a choice!
 - b. This choice is part of your model.
 - c. The confusion matrix is based on a single threshold value.

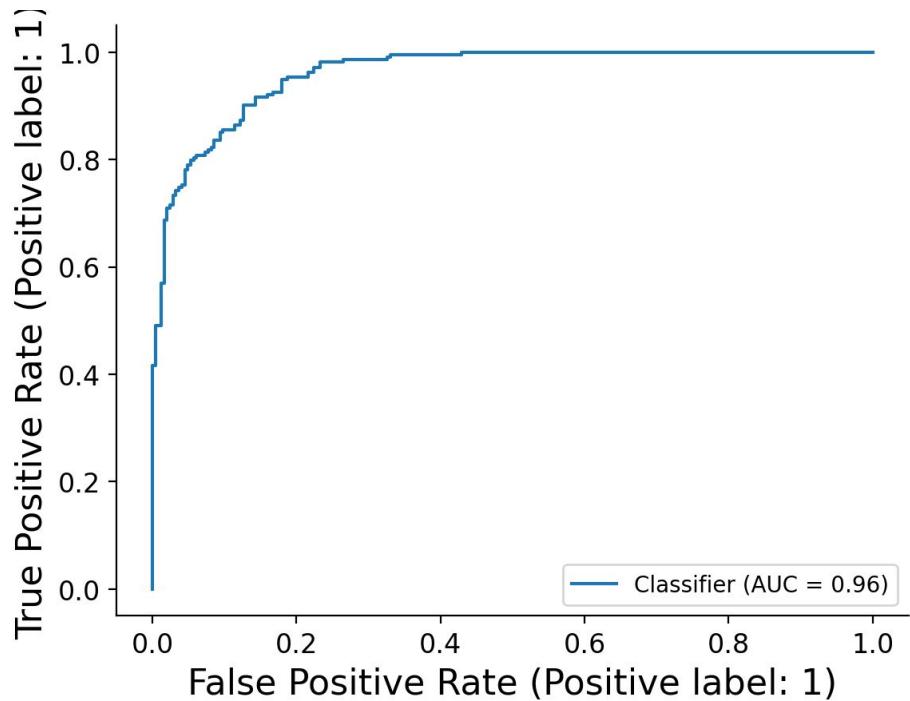
Classification Evaluation - Non-thresholded Measures

- Area Under the Receiver Operating Characteristic (ROC) Curve
 - (Sometimes just called AUC)
- $\text{TPR} = \text{Recall} = \text{TP} / (\text{TP} + \text{FN})$
- $\text{FPR} = \text{TP} / (\text{FP} + \text{TN})$
- Up and to the left is good
- Each point corresponds to a different threshold
- Typically, roughly concave.



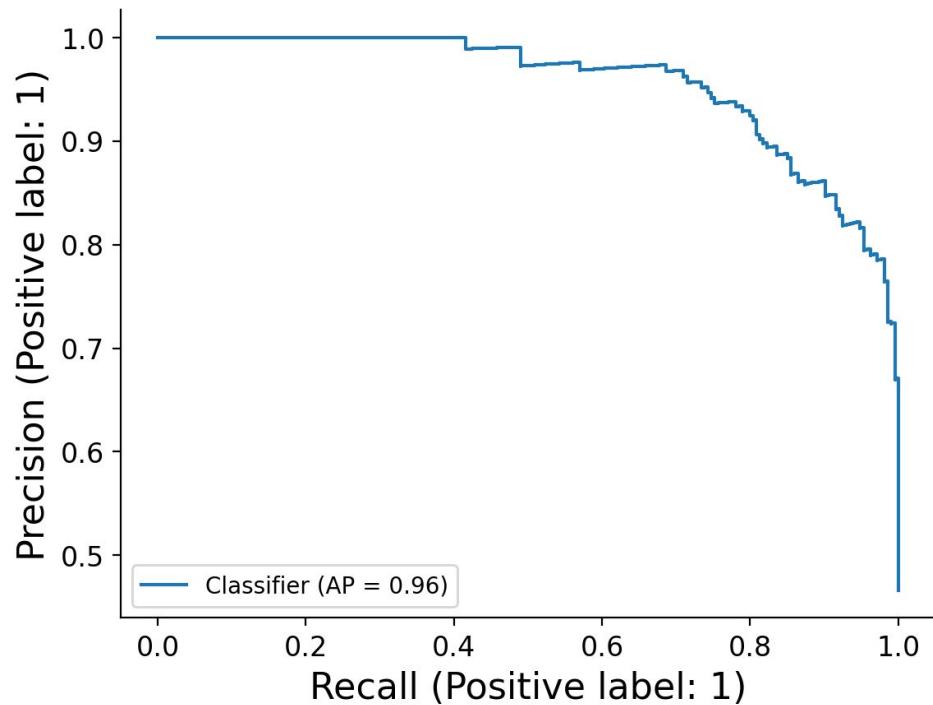
Classification Evaluation - Non-thresholded Measures

- Baseline AUROC is 0.5.
 - i.e. a random classifier
- Good, overall measure of model performance.
- Mitigates much of the impact from imbalance classes.
- Can interpret as “If you randomly draw 1 pos and 1 neg sample, what’s the probability you rank them correctly?”



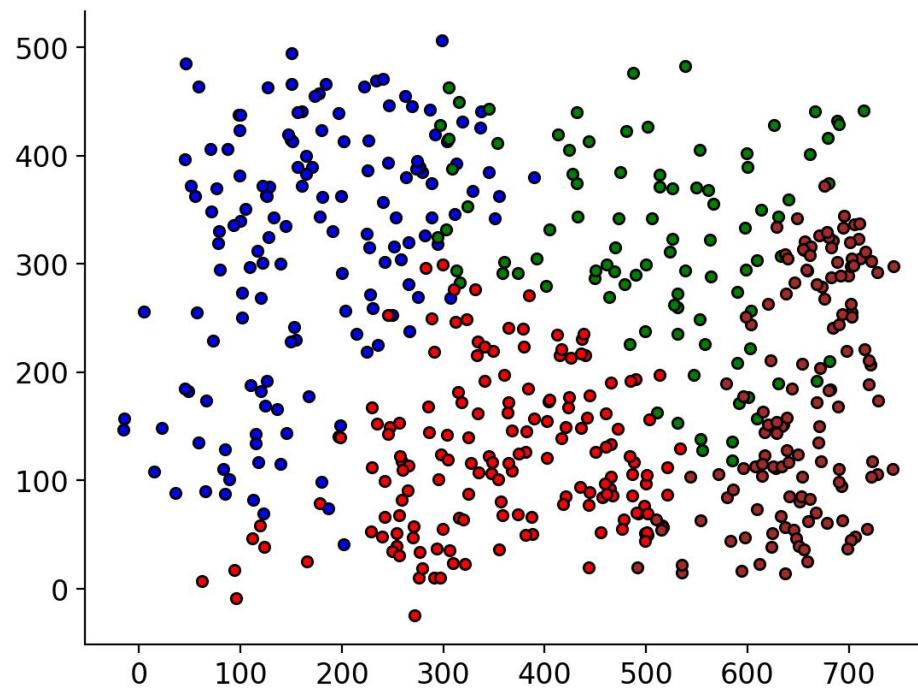
Classification Evaluation - Non-thresholded Measures

- Area Under the Precision Recall Curve (PR curve)
 - “Average Precision”
- Up and to the right is good.
- Poor models exhibit very non-concave PR curves
- Helpful for picking a threshold.

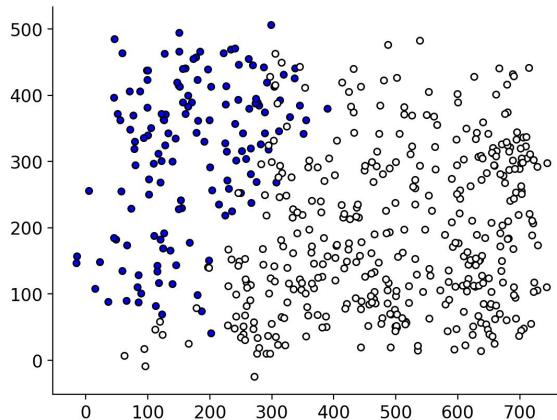


Multiclass Classification

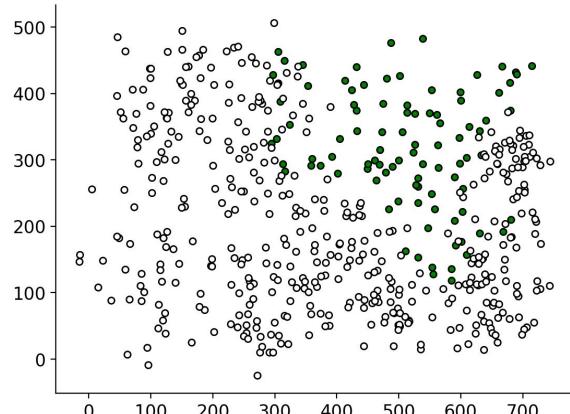
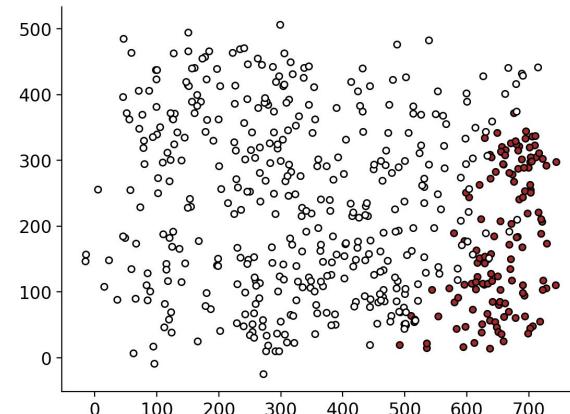
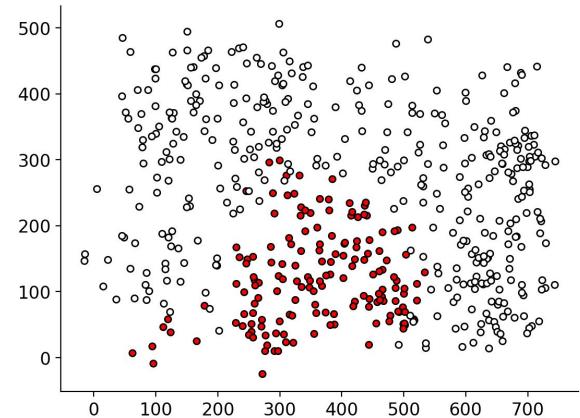
Multiclass Classification



Multiclass Classification - One vs. Rest

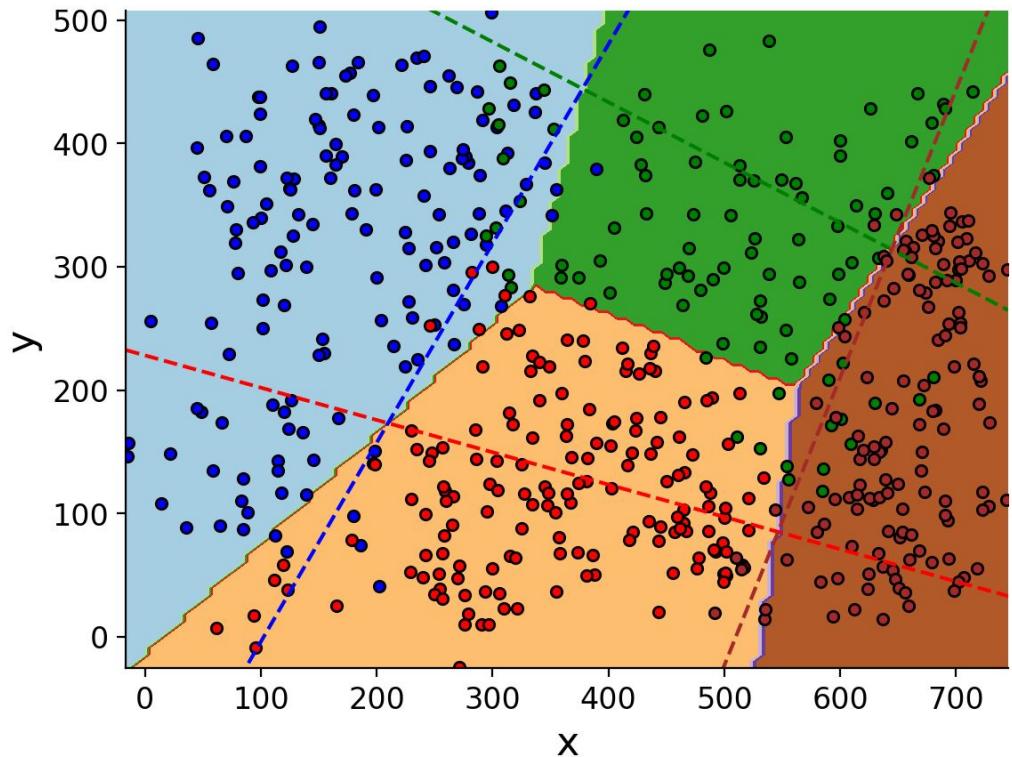


- For N classes, train N binary classifiers.
- Treat each a Nth class as pos label, all other classes as negative label.



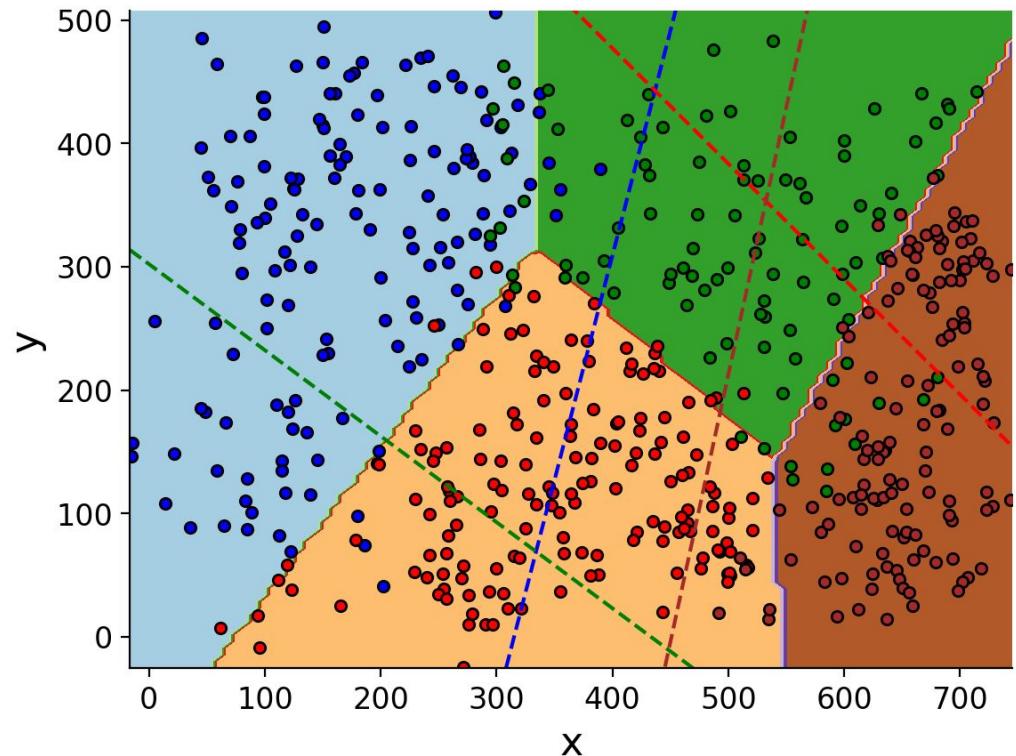
Multiclass Classification - One vs. Rest

- Max classifier score gives the predicted label.



Multiclass Classification - Multinomial

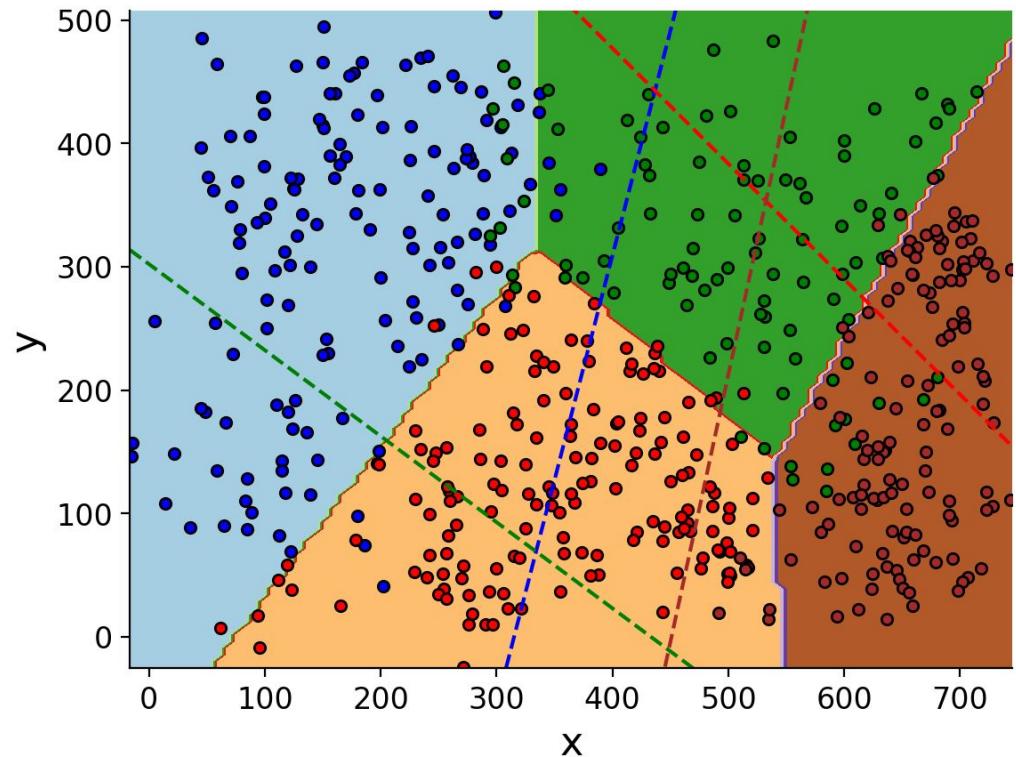
- AKA cross entropy loss



Multiclass Classification - Multinomial / Cross Entropy

$$\mathcal{L}_{CE} = \sum_{i=1}^N \left(- \sum_{c=1}^C y_{ic} \log \left(s \left(\vec{x}_i \cdot \vec{\beta}_c \right) \right) \right)$$

$$s \left(\vec{x}_i \cdot \vec{\beta}_{c'} \right) = \frac{e^{\vec{x}_i \cdot \vec{\beta}_{c'}}}{\sum_{c=1}^C e^{\vec{x}_i \cdot \vec{\beta}_c}}$$



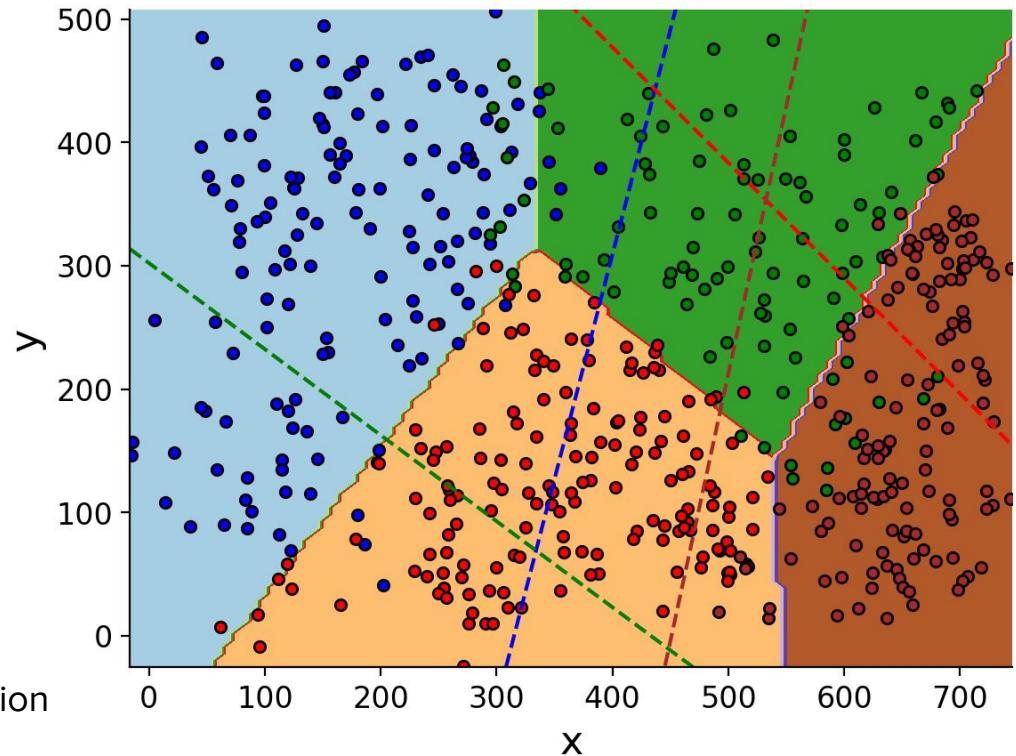
Multiclass Classification - Multinomial / Cross Entropy

$$\mathcal{L}_{CE} = \sum_{i=1}^N \left(- \sum_{c=1}^C y_{ic} \log \left(s \left(\vec{x}_i \cdot \vec{\beta}_c \right) \right) \right)$$

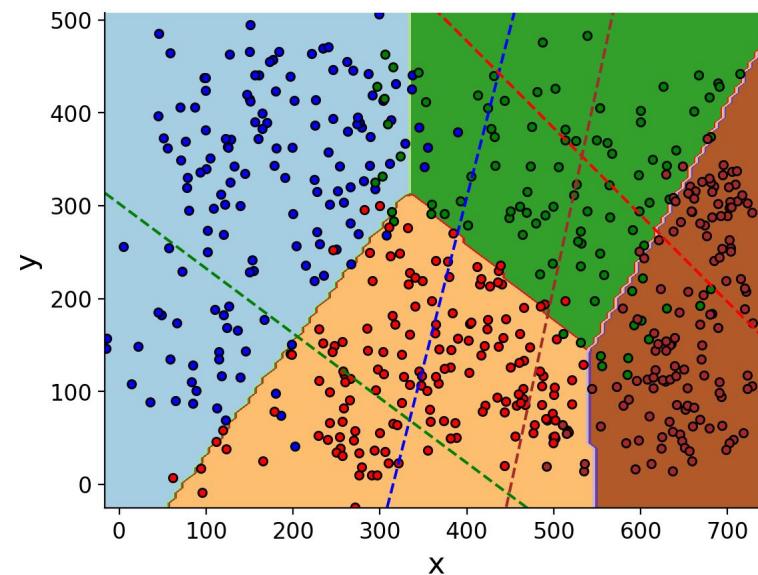
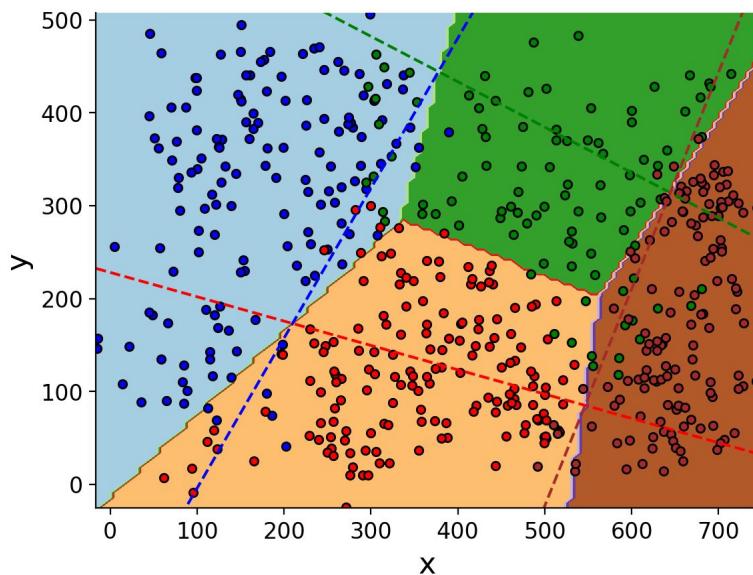
$$s \left(\vec{x}_i \cdot \vec{\beta}_c \right) = \frac{e^{\vec{x}_i \cdot \vec{\beta}_c}}{\sum_{c=1}^C e^{\vec{x}_i \cdot \vec{\beta}_c}}$$

~ Probability sample
belongs to class c'

Softmax function



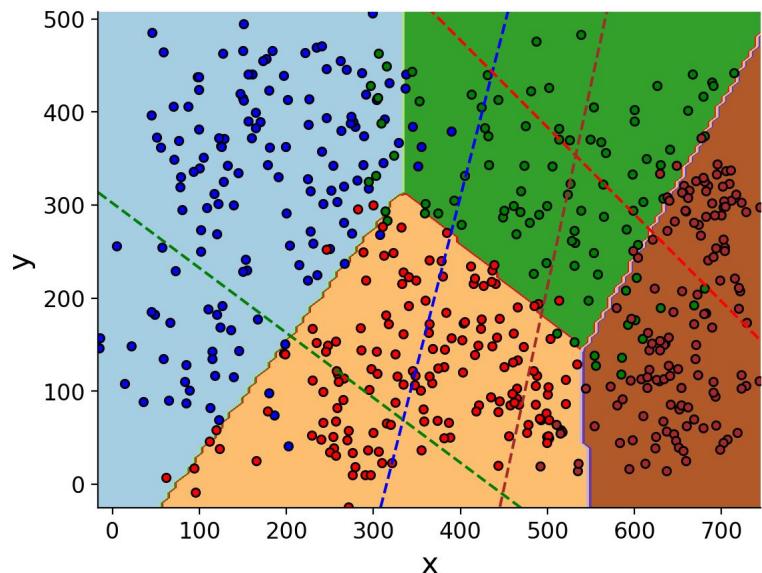
Multiclass Classification - OVR vs Multinomial



Multiclass Classification Evaluation

Confusion Matrix

Predicted Labels \ True Labels	0	1	2	3
0	135	6	6	0
1	9	71	3	13
2	6	5	150	0
3	0	5	5	121



Multiclass Classification Evaluation

Accuracy is now harder!

Binary: For 50/50 labels, random classifier gets 50% accuracy

Multiclass: for equal numbers of C labels, random classifier gets $1 / C$ accuracy

For 4 classes, random classifier gets 25% accuracy.

Multiclass Classification Evaluation

- Precision, recall, and F1 all have multiclass equivalents.
- You can compute them for any individual class.

$$Precision_c = \frac{TP_c}{TP_c + FP_c}$$

- For computing across all classes:
 - Macro-average: Calculate metric for each class and then take the average across all classes.

$$Precision = \frac{1}{C} \sum_c^C \frac{TP_c}{TP_c + FP_c}$$

- Micro-average: Compute numerators and denominators for all classes.

$$Precision = \frac{\sum_c TP_c}{\sum_c TP_c + FP_c}$$

Language Models are Multiclass Classifiers

- Starting Data: “Sounds like fun!”

Language Models are Multiclass Classifiers

- Starting Data: “Sounds like fun!”
- Break it up into data + answers where we try to predict the next word:

Data	Answer
Sounds	like
Sounds like	fun
Sounds like fun	!

Language Models are Multiclass Classifiers

- Starting Data: “Sounds like fun!”
- Break it up into data + answers where we try to predict the next word:

Data	Answer
Sounds	like
Sounds like	fun
Sounds like fun	!

- One-hot-encode words into numbers:

0	0	1	0	1	1	0	1	...
hi	you	like	dog	fun	sounds	me	!	

Language Models are Multiclass Classifiers

- Starting Data: “Sounds like fun!”
- Break it up into data + answers where we try to predict the next word:

Data	Answer
Sounds	like
Sounds like	fun
Sounds like fun	!

- One-hot-encode words into numbers:

0	0	1	0	1	1	0	1	...
hi	you	like	dog	fun	sounds	me	!	

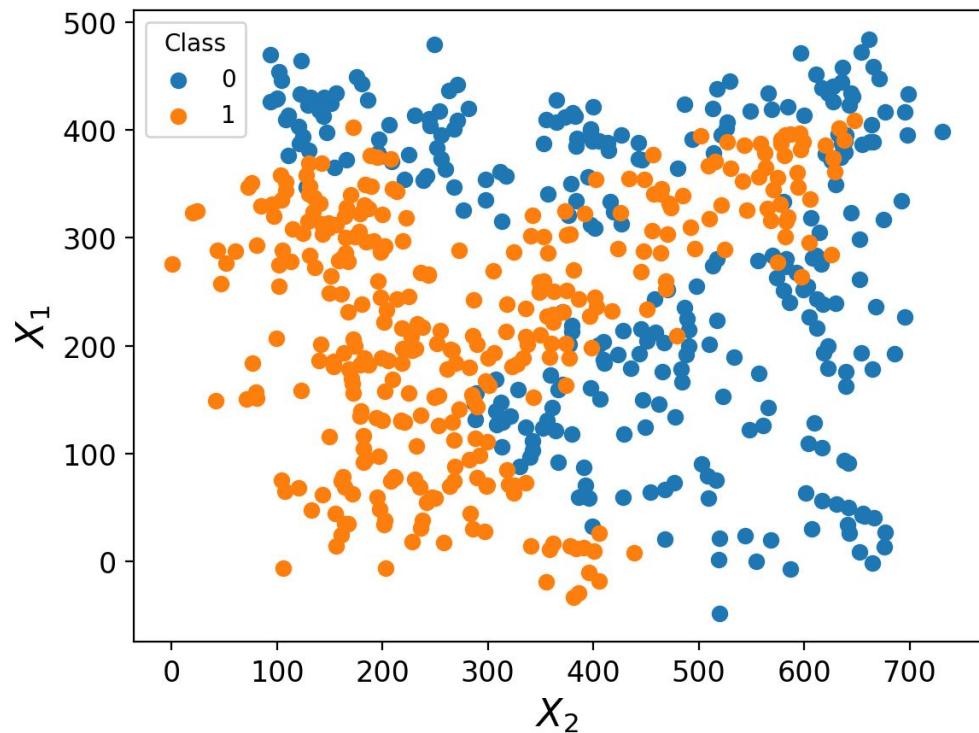
- Feed into a model, treat all words as classes, and predict the answer word (class)

Trees and Ensemble Models

Machine Learning in Financial Engineering
Ethan Rosenthal

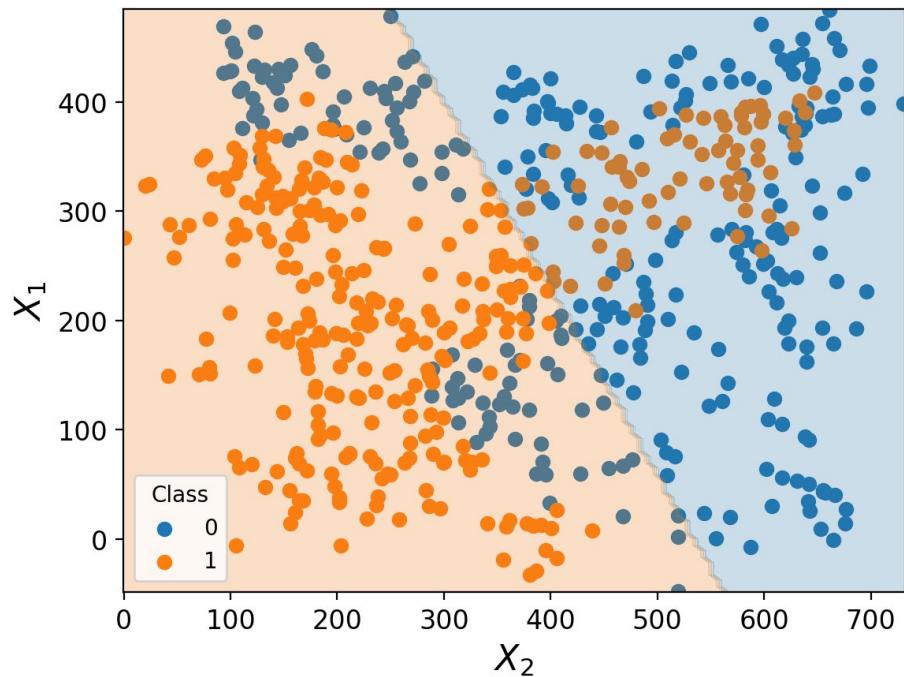
Decision Trees

Limits of Linear Classification



Limits of Linear Classification

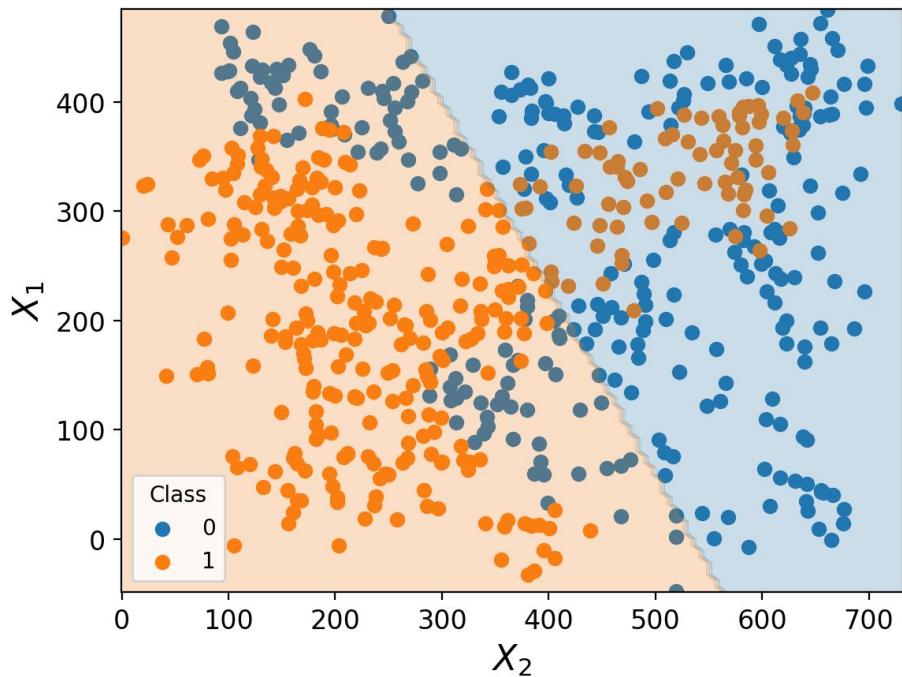
Linear models make linear decision boundaries



Limits of Linear Classification

Linear models make linear decision boundaries

But what if we combined lots of them together?

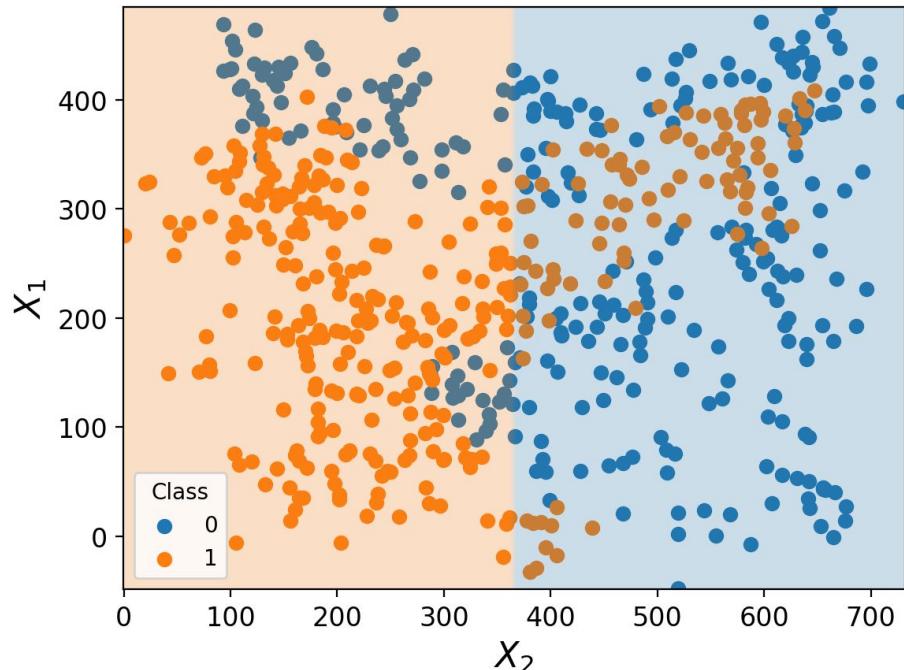
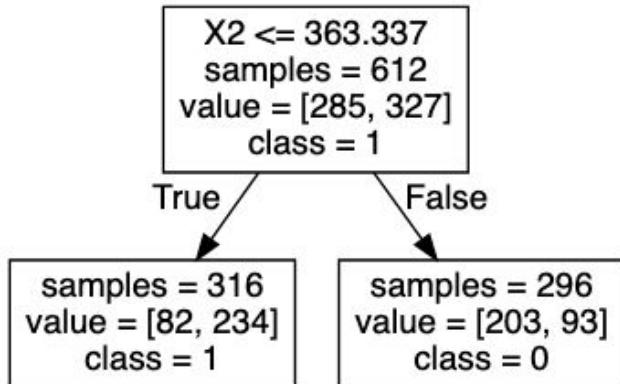


Decision Trees

Start with a very simple “rule”:

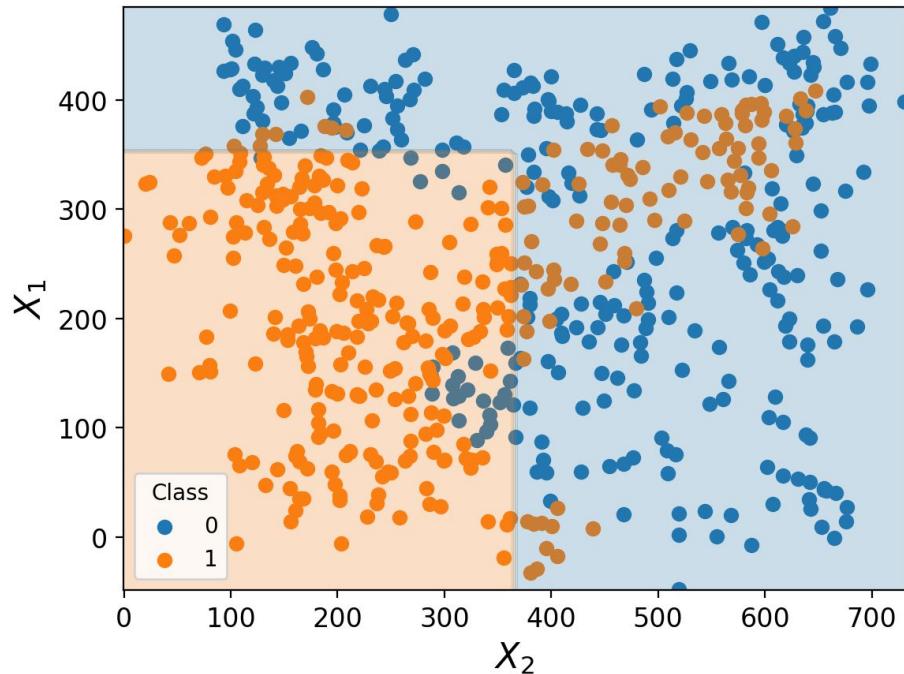
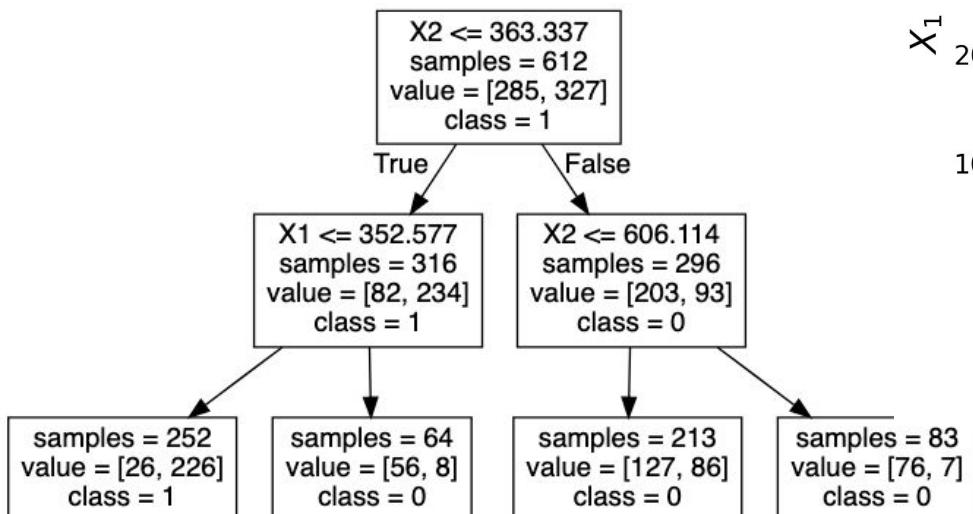
If $X_2 \leq 363.337$,

then Class 1, else Class 0



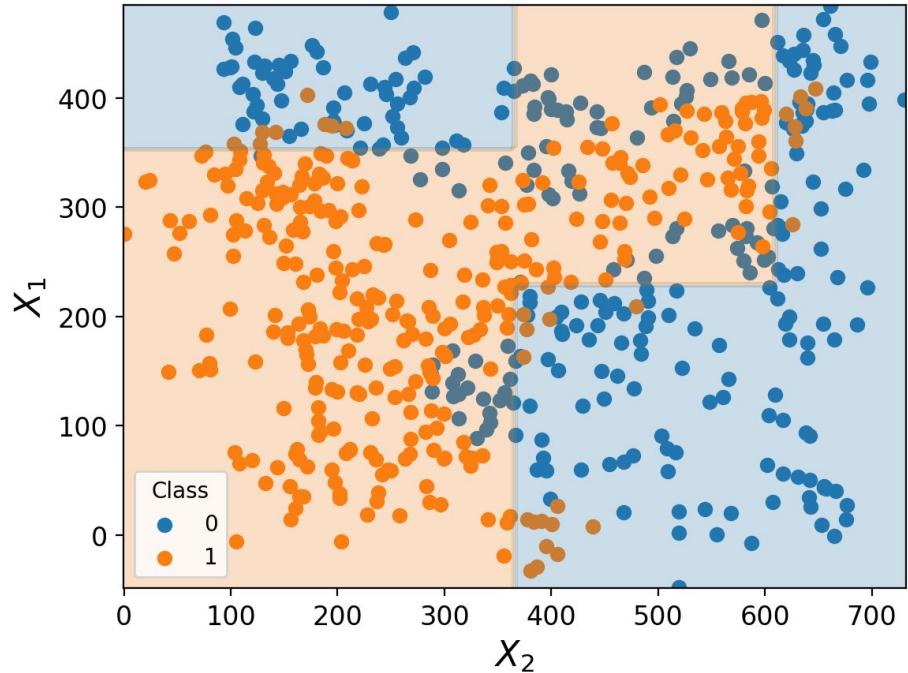
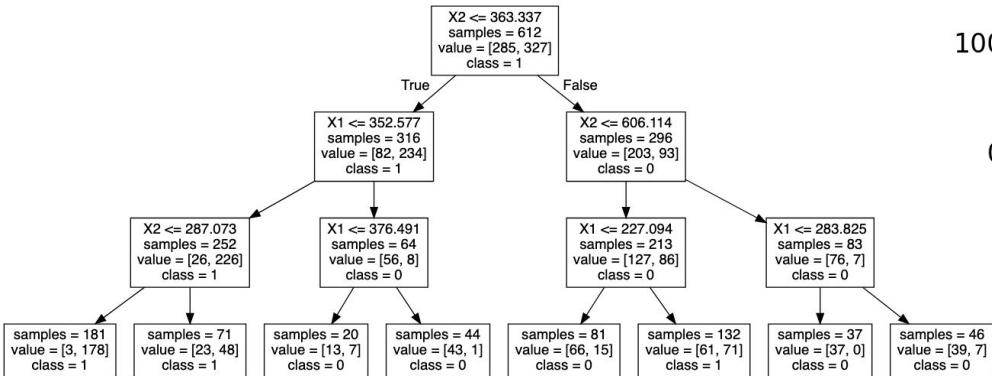
Decision Trees

Add more if/else statements, get more sides to the decision boundary.



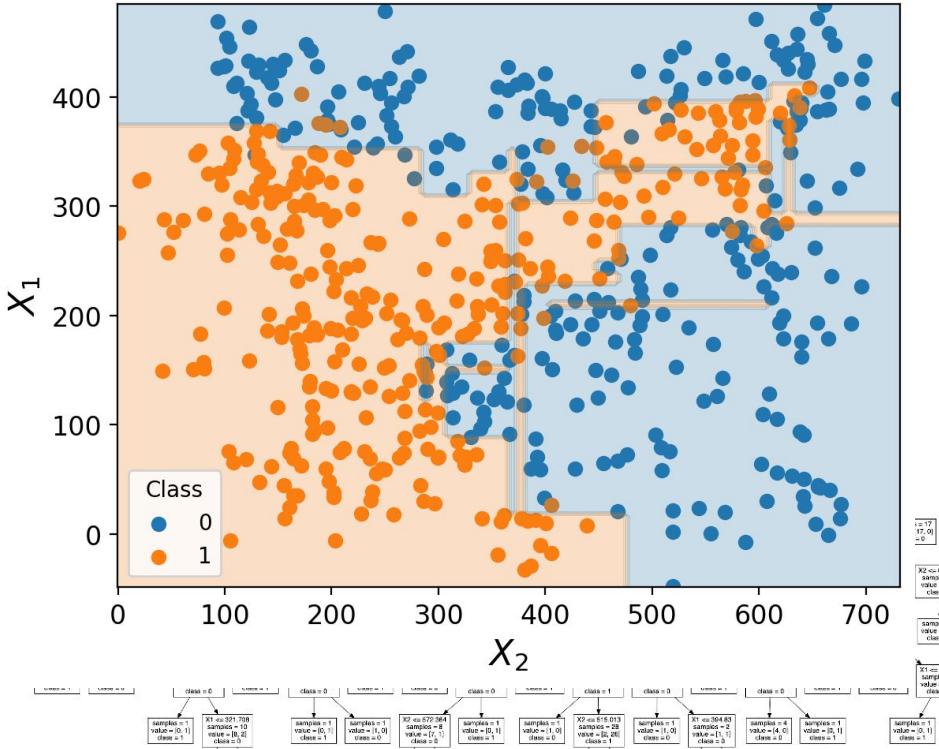
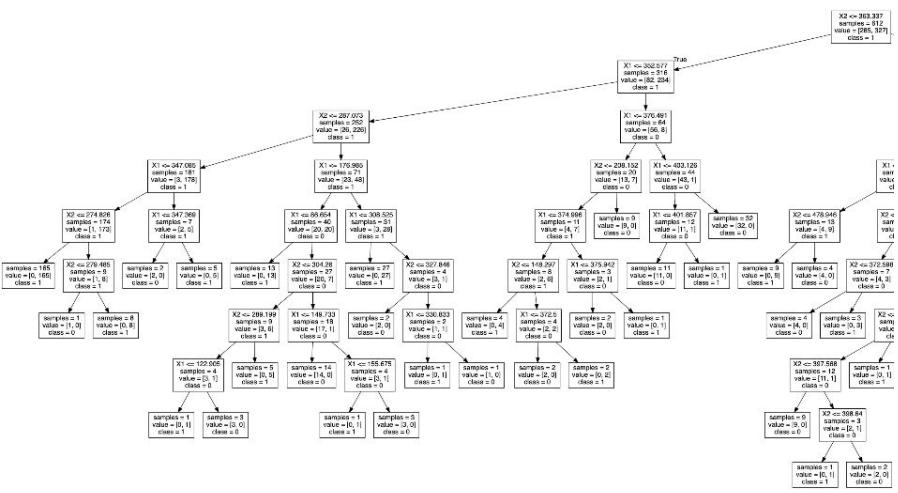
Decision Trees

Add more if/else statements, get more sides to the decision boundary.

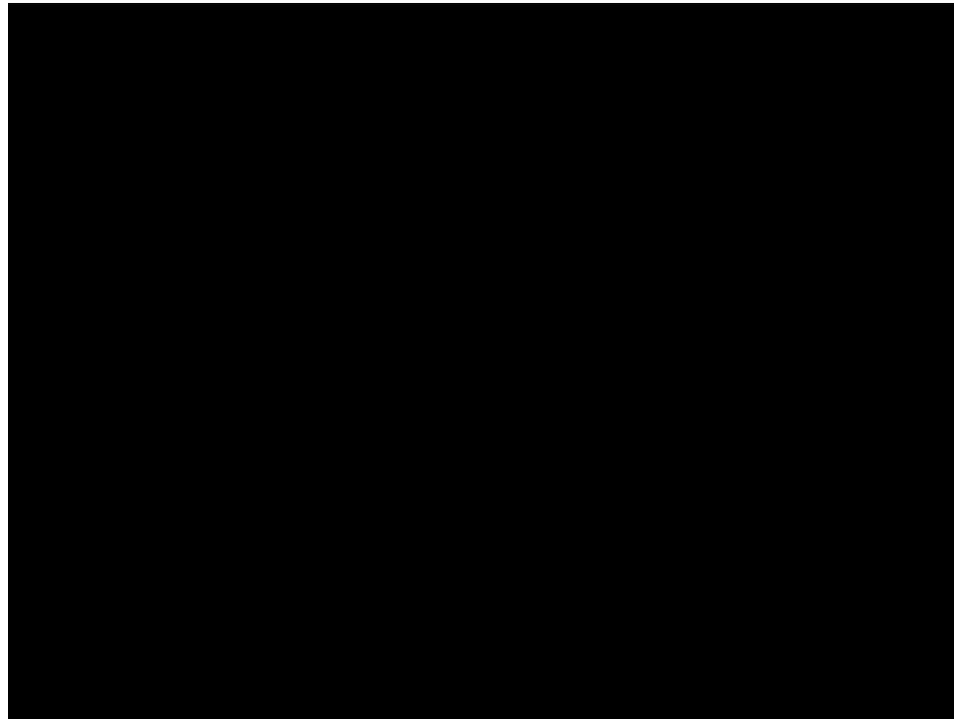


Decision Trees

Add more if/else statements, get more sides to the decision boundary.

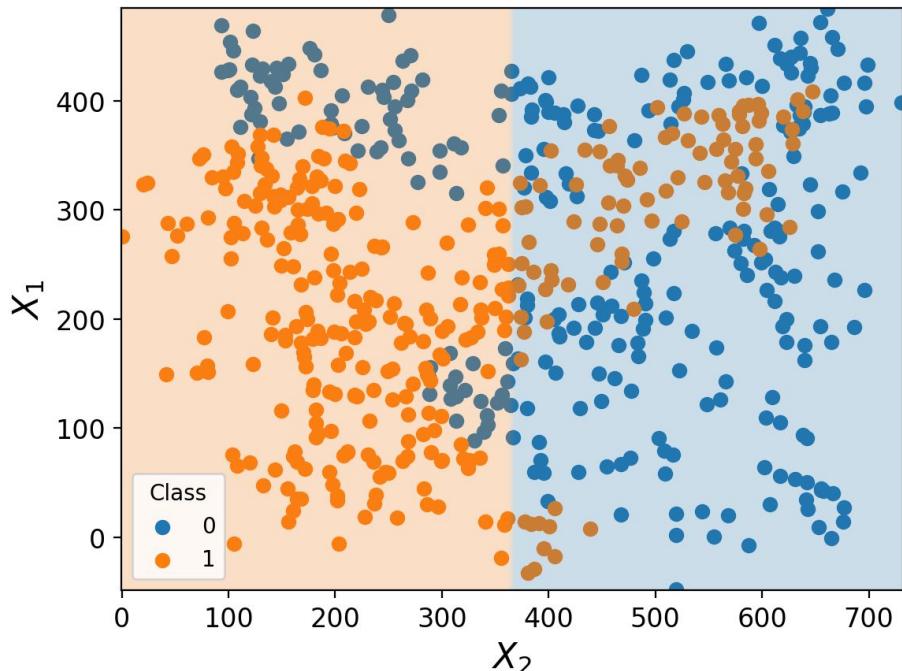
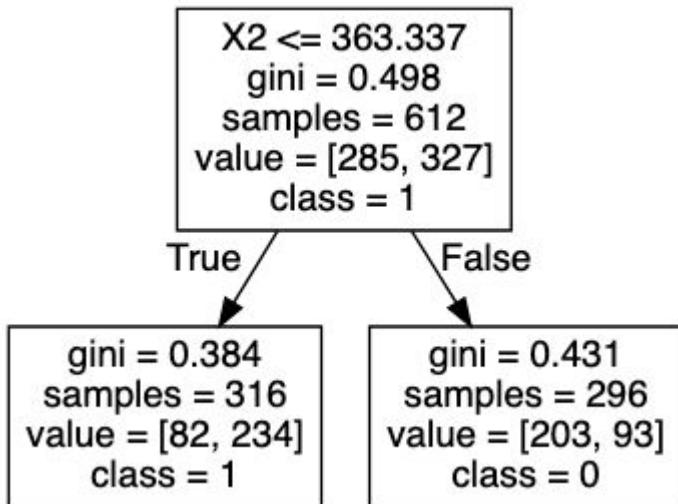


Decision Trees



Decision Trees - How do they Grow?

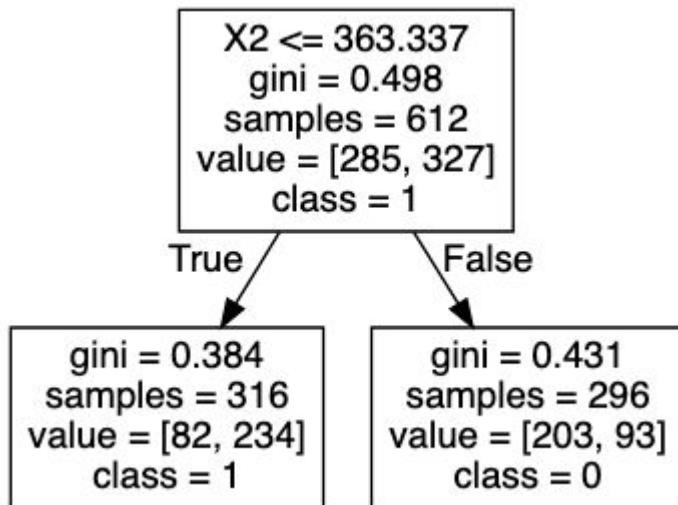
At each node, goal is to find the feature and threshold that maximally splits the classes.



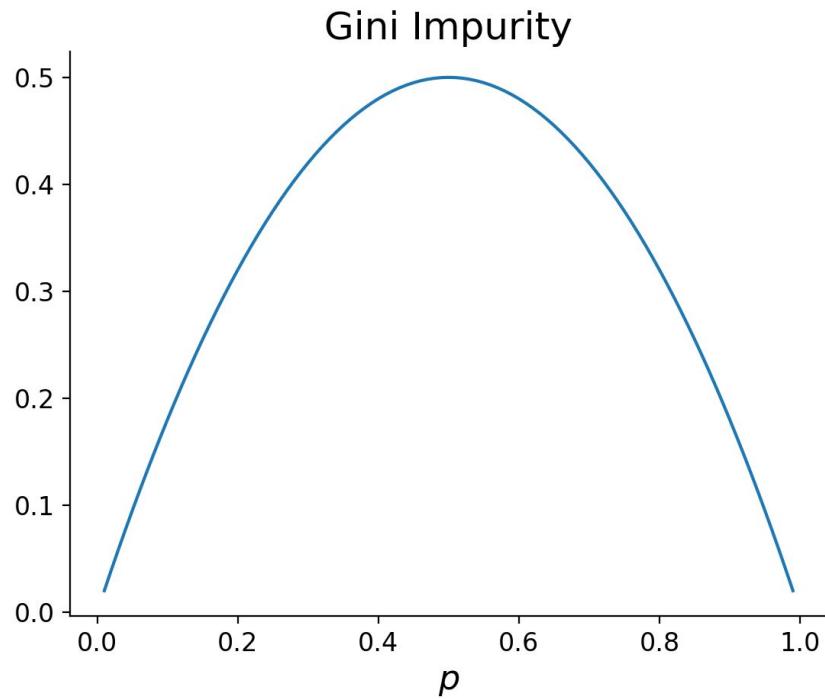
Proportion of samples in the positive class after the node

Decision Trees - How do they Grow?

The Gini Impurity is one measure of how well split the classes are.



$$Gini = 2p(1 - p)$$



Decision Trees - How do they Grow?

Various other
“impurity” measures

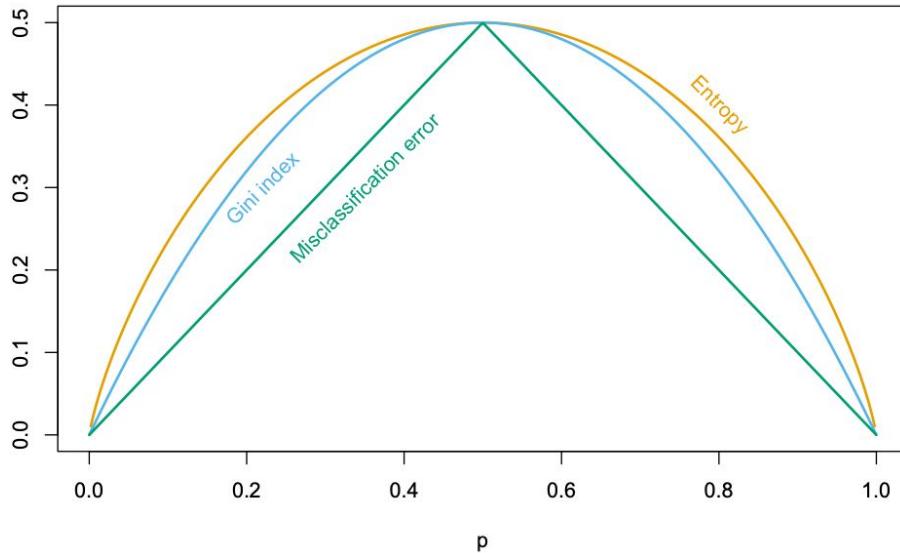
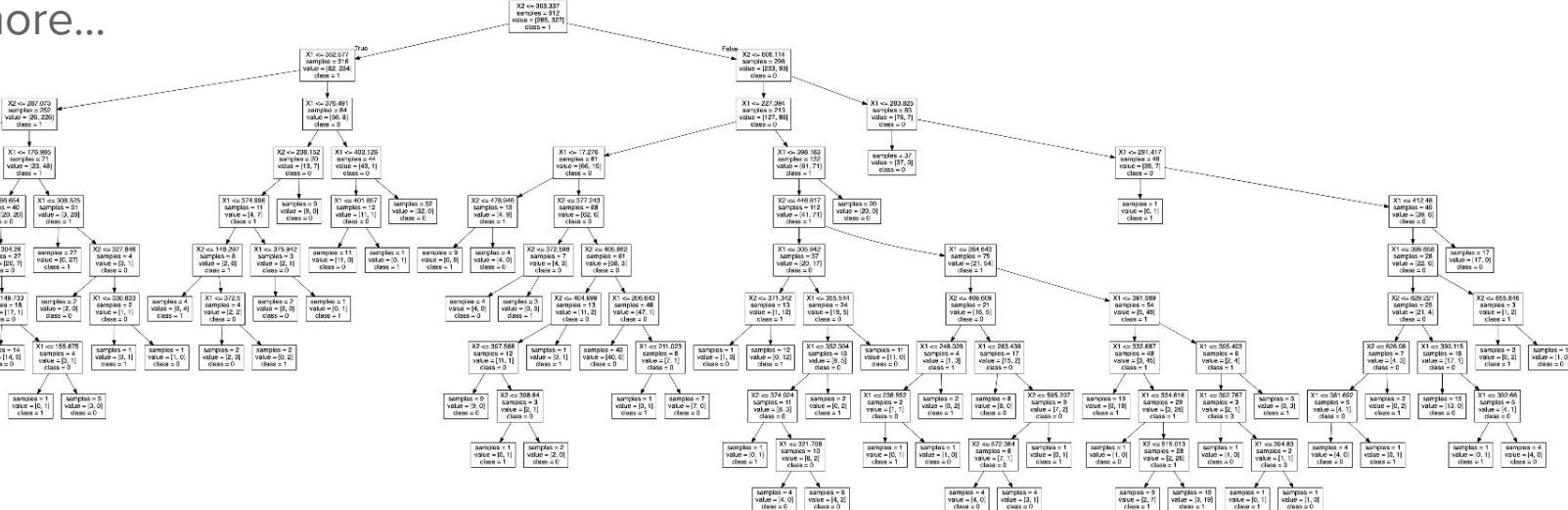


FIGURE 9.3. Node impurity measures for two-class classification, as a function of the proportion p in class 2. Cross-entropy has been scaled to pass through $(0.5, 0.5)$.

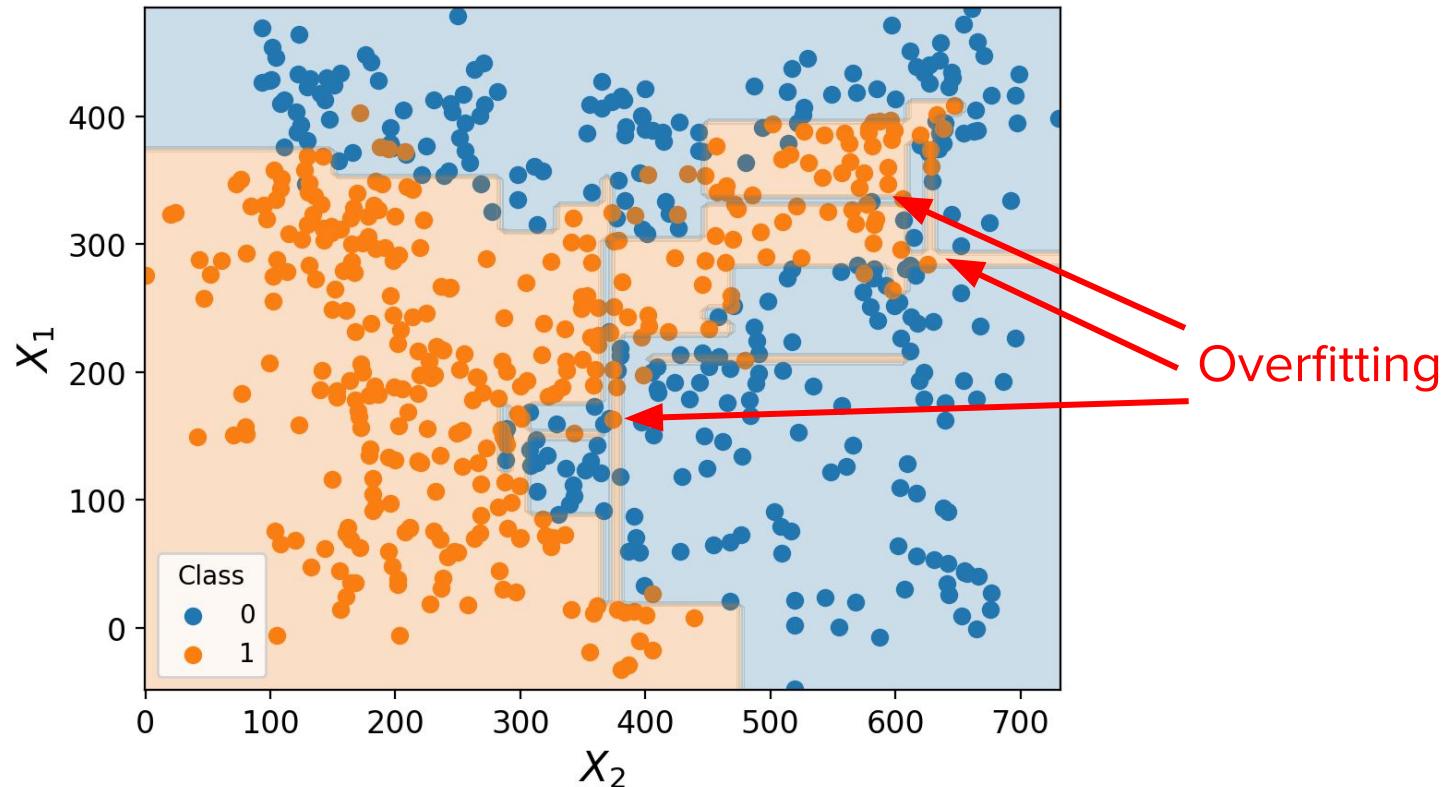
Decision Trees - How do they Grow?

Keep growing the tree until some cutoff criteria:

- Max depth
- Min samples in leaf nodes reached
- Min impurity decrease
- And more...



Limits of Decision Trees



Random Forests

Random Forests - Seeing the Forest for the Trees

Instead of a single decision tree, create many trees (a forest!).

But, induce randomness.

For each tree:

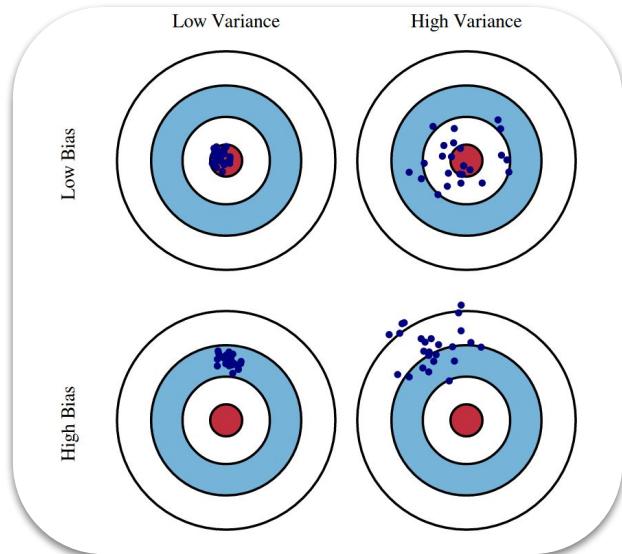
- Generate a *bootstrap* sample of the dataset (i.e. sample with replacement).
- For each node, only consider a subset of the features when deciding what feature to split on.
- The prediction score for each class is the fraction of trees that classify the sample into that class.

Random Forests - Seeing the Forest for the Trees

- Each tree is not as good at predicting as a single decision tree due induced randomness.
- But, the forest helps prevent overfitting.
- This overfitting prevention is often more powerful than the weakness of each tree, leading to a better overall model.
- This is a manifestation of the bias-variance tradeoff.

Random Forests - Seeing the Forest for the Trees

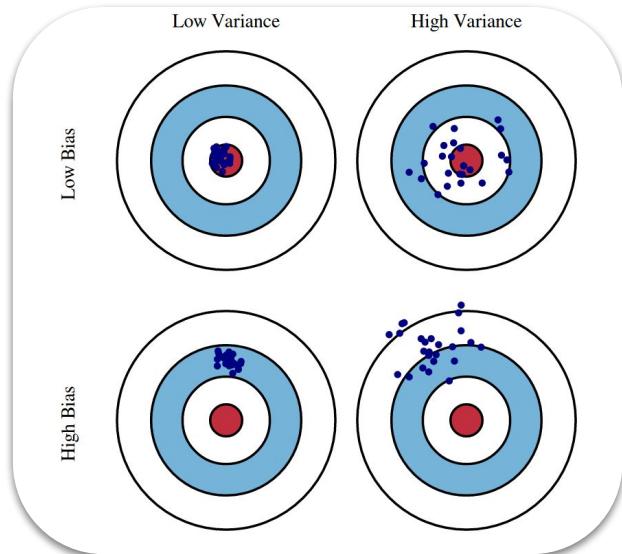
- Each prediction error is a result of:
 - Bias Error: due to our assumptions about the target function.
 - Variance Error: due to the specifics of the dataset.
 - Irreducible Error: nothing we can do here!
- Examples:
 - Regression has low/high bias but a low/high variance.
 - Decision trees have low/high bias but a low/high variance.
 - Random forests have low/high bias but a low/high variance.



The typical bias/variance image in all blog posts on the web!

Random Forests - Seeing the Forest for the Trees

- Each prediction error is a result of:
 - Bias Error: due to our assumptions about the target function.
 - Variance Error: due to the specifics of the dataset.
 - Irreducible Error: nothing we can do here!
- Examples:
 - Regression has low/**high** bias but a **low/high** variance.
 - Decision trees have **low/high** bias but a **low/high** variance (**overfitting**).
 - Random forests have **low/high** bias but a **low/high** variance (less than single trees).



The typical bias/variance image in all blog posts on the web!

Random Forests - Why use them?

- Naturally handle nonlinear relationships in the data.
- Quick to fit.
- Robust (but not immune) to overfitting.
- You don't have to scale your data.
- They can be (kind of) interpretable.
 - See [feature importances](#) which measures the total Gini reduction brought by each feature.
- They just work really well!