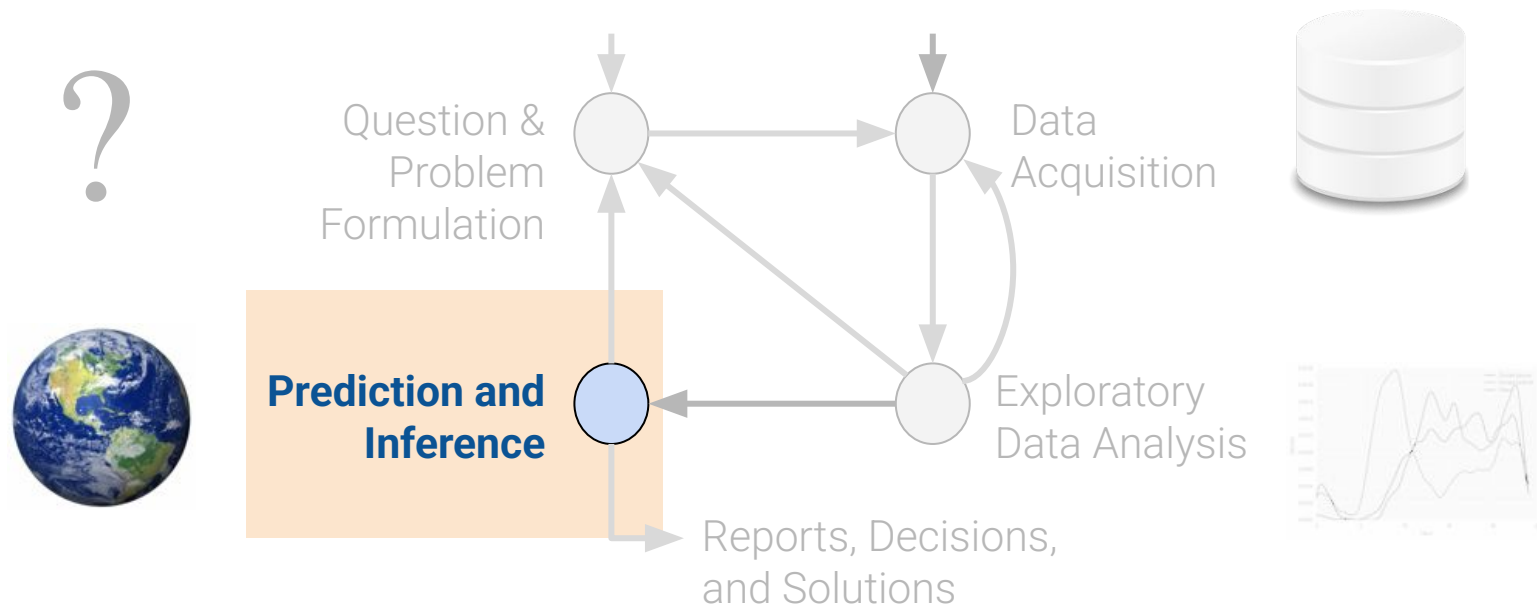# Logistic Regression II

Model Performance.

**Data 100/Data 200, Spring 2022 @ UC Berkeley**

Josh Hug and Lisa Yan

# More Logistic Regression



**(today)**

| Logistic Regression I: | Logistic Regression II: |
|---|---|
| The Model | Linear Separability |
| Cross-Entropy Loss | Accuracy, Precision, Recall |
| The Probabilistic View | Classification Thresholds |

# Today's Roadmap

Lecture 22, Data 100 Spring 2022

Logistic Regression Model, continued

- sklearn demo
- Maximum Likelihood Estimation: high-level (live), detailed (recorded)

Linear separability and Regularization

Performance Metrics

- Accuracy
- Imbalanced Data, Precision, Recall

Adjusting the Classification Threshold

- A case study
- ROC curves, and AUC

[Extra] Detailed MLE, Gradient Descent, PR curves

# Logistic Regression Model, continued

**Logistic Regression Model, continued**

- sklearn demo
- Maximum Likelihood Estimation: high-level (live), detailed (recorded)

Linear separability and Regularization

Performance Metrics

- Accuracy
- Imbalanced Data, Precision, Recall

Adjusting the Classification Threshold

- A case study
- ROC curves, and AUC

[Extra] Detailed MLE, Gradient Descent, PR curves

Lecture 22, Data 100 Spring 2022

```
from sklearn.linear_model import LogisticRegression
model = LogisticRegression(fit_intercept=False)
model.fit(X, Y)
```

**Demo**

| Task/Model |
|---|

Binary Classification ($y \in \{0, 1\}$)

$$\hat{P}_\theta(Y = 1|x) = \sigma(x^T\theta)$$

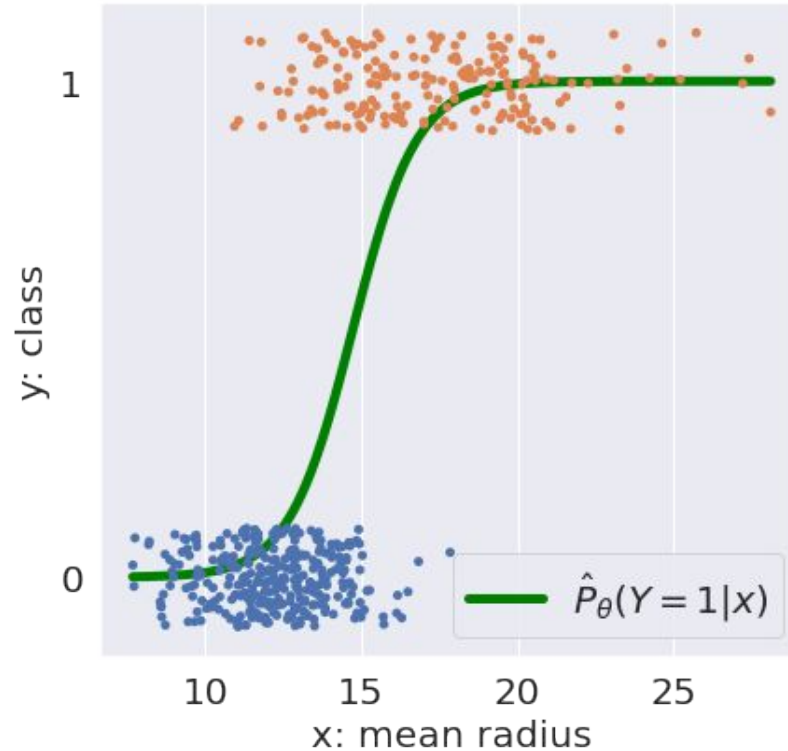| Fit to objective function |
|---|

Average Cross-Entropy Loss

$$-\frac{1}{n}\sum_{i=1}^{n}\left(y_i\log(\sigma(X_i^T\theta) + (1 - y_i)\log(1 - \sigma(X_i^T\theta)))\right)$$

+ regularization

For logistic regression, sklearn applies regularization by default. We'll see why soon.

5

# Demo

## Sklearn: Predict Probabilities

```
model.predict_proba(X) # probs for all classes
model.classes_         # array([0, 1])
```

```
model.predict_proba(X) # probs for all classes
model.classes_         # array([0, 1])
```

```
model.predict(X)              # predict 1 or 0
```

$$\hat{y} = \text{classify}(x) = \begin{cases} 1 & \hat{P}_\theta(Y = 1|x) \geq 0.5 \\ 0 & \text{otherwise} \end{cases}$$

Equivalent "otherwise" condition: $\hat{P}_\theta(Y = 0|x) \geq 0.5$

**Demo**

Interpret:  Given the input feature x:
If Y is more likely to be 1 than 0,
then predict $\hat{y} = 1$.
Else predict 0.

|   | X | Y | P(Y = 1 \| x) | Y_hat |
|---|---|---|---|---|
| 0 | 25.220 | 1 | 0.999965 | 1 |
| 1 | 13.480 | 1 | 0.226448 | 0 |
| 2 | 11.290 | 0 | 0.033174 | 0 |
| 3 | 12.860 | 0 | 0.137598 | 0 |
| 4 | 19.690 | 1 | 0.992236 | 1 |

7

# [High-Level] Maximum Likelihood Estimation

**Minimizing cross-entropy loss** is equivalent to **maximizing the likelihood of the training data**.

Assumption: all data are independent Bernoulli random variables.

$$\underset{\theta}{\textbf{argmin}} \quad -\frac{1}{n}\sum_{i=1}^{n}\left(y_i \log(\sigma(X_i^T\theta) + (1-y_i)\log(1-\sigma(X_i^T\theta))\right)$$
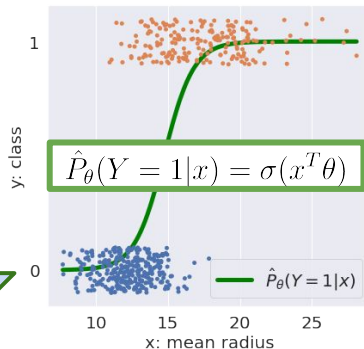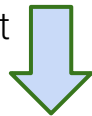
For logistic regression, let $p_i = \sigma(X_i^T\theta)$

$$\underset{p_1, p_2, \ldots, p_n}{\textbf{argmax}} \quad \prod_{i=1}^{n} \underbrace{p_i^{y_i}(1-p_i)^{(1-y_i)}}_{\text{Prob. that i-th response is yi}}$$

**Main takeaway**: The optimal theta that minimizes mean cross-entropy loss "pushes" all probabilities in the direction of the true class.



Want $\sigma(x^T\theta) \to 1$

$\hat{P}_\theta(Y=1|x) = \sigma(x^T\theta)$

Want $\sigma(x^T\theta) \to 0$

# [High-Level] Maximum Likelihood Estimation

**Minimizing cross-entropy loss** is equivalent to **maximizing the likelihood of the training data**.

$$\underset{\theta}{\textbf{argmin}} \quad -\frac{1}{n}\sum_{i=1}^{n}\left(y_i \log(\sigma(X_i^T\theta) + (1-y_i)\log(1-\sigma(X_i^T\theta)))\right)$$

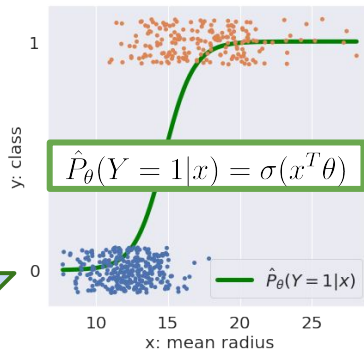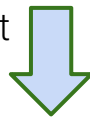For logistic regression, let $p_i = \sigma(X_i^T\theta)$

Assumption: all data are independent Bernoulli random variables

$$\underset{p_1, p_2, \ldots, p_n}{\textbf{argmax}} \quad \prod_{i=1}^{n} \underbrace{p_i^{y_i}(1-p_i)^{(1-y_i)}}_{\text{Prob. that i-th response is yi}}$$

**Main takeaway**: The optimal theta that minimizes mean cross-entropy loss "pushes" all probabilities in the direction of the true class.



Want $\sigma(x^T\theta) \to 1$

$\hat{P}_\theta(Y=1|x) = \sigma(x^T\theta)$

Want $\sigma(x^T\theta) \to 0$

# Linear separability and Regularization

Lecture 22, Data 100 Spring 2022

**Demo**

```
from sklearn.linear_model import LogisticRegression
model = LogisticRegression(fit_intercept=False)
model.fit(X, Y)
```

Task/Model

Binary Classification ($y \in \{0, 1\}$)

$$\hat{P}_\theta(Y = 1|x) = \sigma(x^T\theta)$$

Fit to objective function

Average Cross-Entropy Loss

$$-\frac{1}{n}\sum_{i=1}^{n}\left(y_i\log(\sigma(X_i^T\theta) + (1 - y_i)\log(1 - \sigma(X_i^T\theta)))\right)$$

+ **regularization**

Why does sklearn always apply regularization?

11

# Linear Separability

A classification dataset is said to be **linearly separable** if there exists a hyperplane **among input features x** that separates the two classes y.

If there is one feature; the input feature is **1-D**.
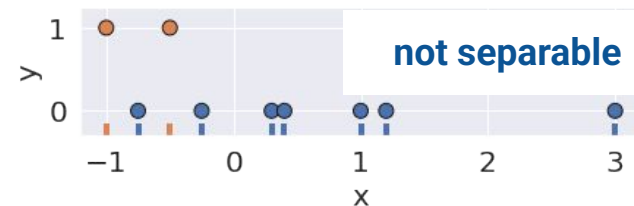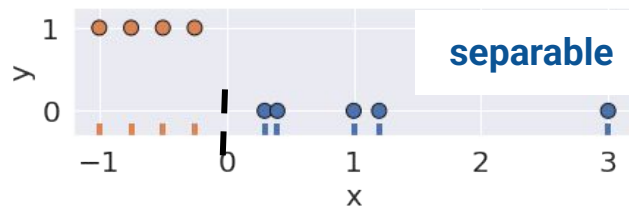- Class label is not a feature; it is output.
- Use rug plot to see separability.

# Linear Separability

A classification dataset is said to be **linearly separable** if there exists a hyperplane **among input features x** that separates the two classes y.

If there is one feature; the input feature is **1-D**.
- Class label is not a feature; it is output.
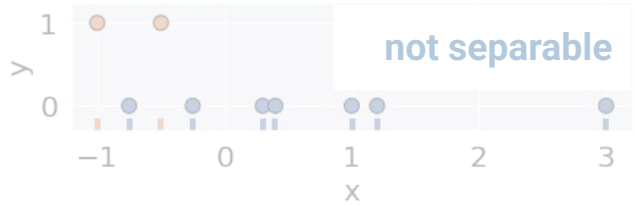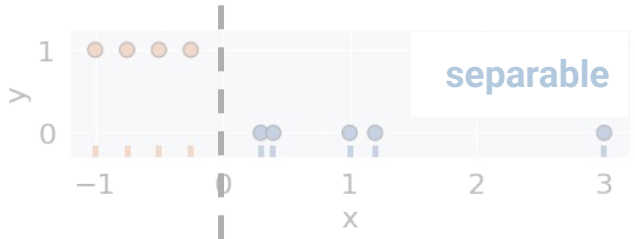- Use rug plot to see separability.



If there are two features, the input feature is **2-D**. Use scatter plot to see separability.

13

# Linearly Separability Creates Diverging Weights

Consider the simplified logistic regression model fit to the toy data:

$$\hat{P}_\theta(Y = 1|x) = \sigma(\theta x) = \frac{1}{1 + e^{-\theta x}}$$

What will be the optimal weight theta? Why?

A. $\hat{\theta} = -1$    C. $\hat{\theta} \to -\infty$

B. $\hat{\theta} = 1$    D. $\hat{\theta} \to \infty$

[Hint] The optimal theta should "push" probabilities in the direction of the true class:

$\bullet$  $\hat{P}_\theta(Y = 1|x = -1) = \dfrac{1}{1 + e^\theta}$   $\to 1$

$\bullet$  $\hat{P}_\theta(Y = 1|x = 1) = \dfrac{1}{1 + e^{-\theta}}$   $\to 0$

The Data

(-1,1)

(1,0)

| x | y |
|---|---|
| -1 | 1 |
| 1 | 0 |

# Linearly Separability Creates Diverging Weights

Consider the simplified logistic regression model fit to the toy data:

$$\hat{P}_\theta(Y = 1|x) = \sigma(\theta x) = \frac{1}{1 + e^{-\theta x}}$$

What will be the optimal weight theta? Why?

A. $\hat{\theta} = -1$     C. $\hat{\theta} \to -\infty$

B. $\hat{\theta} = 1$     D. $\hat{\theta} \to \infty$

[Hint] The optimal theta should "push" probabilities in the direction of the true class:

⬤ $\hat{P}_\theta(Y = 1|x = -1) = \dfrac{1}{1 + e^{\theta}}$    $\to 1$

⬤ $\hat{P}_\theta(Y = 1|x = 1) = \dfrac{1}{1 + e^{-\theta}}$    $\to 0$

happens as $\hat{\theta} \to -\infty$

The Data

| x | y |
|---|---|
| -1 | 1 |
| 1 | 0 |

(-1,1)

(1,0)

1

-1      0      1

# Linearly Separability Creates Diverging Weights

Consider the simplified logistic regression model fit to the toy data:

$$\hat{P}_\theta(Y = 1|x) = \sigma(\theta x) = \frac{1}{1 + e^{-\theta x}}$$

What will be the optimal weight theta? Why?
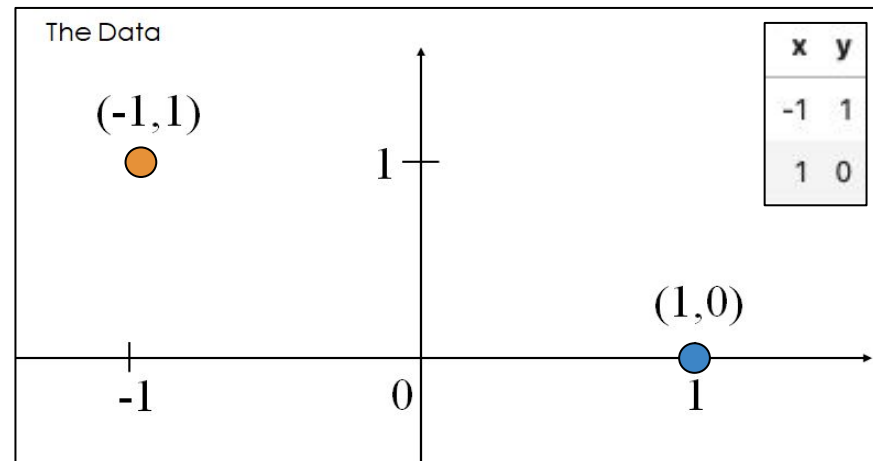
A. $\hat{\theta} = -1$

B. $\hat{\theta} = 1$

C. $\hat{\theta} \to -\infty$

D. $\hat{\theta} \to \infty$

The Data

| x | y |
|---|---|
| -1 | 1 |
| 1 | 0 |

(-1,1)

(1,0)

1

-1    0    1

[Hint] The optimal theta should "push" probabilities in the direction of the true class:

$$\hat{P}_\theta(Y = 1|x = -1) = \frac{1}{1 + e^\theta} \quad \to 1$$

$$\hat{P}_\theta(Y = 1|x = 1) = \frac{1}{1 + e^{-\theta}} \quad \to 0$$

happens as $\hat{\theta} \to -\infty$

Mean Loss on Toy Data

direction of gradient

(Impossible to see, but) plateau is slightly tilted downwards.

Loss approaches 0 as theta decreases.

# Linearly Separability Creates Diverging Weights

Consider the simplified logistic regression model fit to the toy data:

$\hat{\theta} \to -\infty$:

$$\hat{P}_\theta(Y = 1 | x = -1) = \frac{1}{1 + e^\theta} \to 1$$

$$\hat{P}_\theta(Y = 1 | x = 1) = \frac{1}{1 + e^{-\theta}} \to 0$$

$$\hat{P}_\theta(Y = 1 | x) = \sigma(\theta x) \to \begin{cases} 1 & \text{if } x < 0 \\ 0 & \text{if } x \geq 0 \end{cases}$$

The Data



(-1,1)

(1,0)

1

-1    0    1

# Linearly Separability Creates Diverging Weights

Consider the simplified logistic regression model fit to the toy data:

$\hat{\theta} \rightarrow -\infty$:

$$\hat{P}_\theta(Y = 1|x = -1) = \frac{1}{1 + e^\theta} \rightarrow 1$$

$$\hat{P}_\theta(Y = 1|x = 1) = \frac{1}{1 + e^{-\theta}} \rightarrow 0$$



The Data

(-1,1)

(1,0)

$$\hat{P}_\theta(Y = 1|x) = \sigma(\theta x) \rightarrow \begin{cases} 1 & \text{if } x < 0 \\ 0 & \text{if } x \geq 0 \end{cases}$$

This model is **overconfident**.

- Consider a new point (0.5, **1**).
- The model incorrectly says **p = 0**, so it predicts **0**.

**Typo fixed 5/3**



YES... HA HA HA... *YES!*

Diverging weights

$$\overset{=1}{-(y \log(p)} + (1 - y) \log(1 - p))$$

$$\rightarrow 1 \log(0)$$    **Loss is infinite.**

Divergent weights (i.e., | θ | →∞) occur with **linearly separable** data.

"Overconfidence" is a particularly dangerous version of overfitting.

# Regularized Logistic Regression

To avoid large weights (particularly on linearly separable data), use **regularization**.

- As with linear regression, standardize features first.

$$\underset{\theta}{\operatorname{argmin}} -\frac{1}{n} \sum_{i=1}^{n} \left( y_i \log(\sigma(X_i^T \theta) + (1 - y_i) \log(1 - \sigma(X_i^T \theta))) \right) + \lambda \sum_{j=1}^{d} \theta_j^2$$



```
# sklearn defaults
model = LogisticRegression(
         penalty='l2', C=1.0, …)
model.fit()
```

Regularization hyperparameter C is the inverse of $\lambda$. C = 1 / $\lambda$.

Set C big for minimal regularization, e.g., C=300.0.

19

# Performance Metrics

Lecture 22, Data 100 Spring 2022

Logistic Regression Model, continued

- sklearn demo
- Maximum Likelihood Estimation: high-level (live), detailed (recorded)

Linear separability and Regularization

## Performance Metrics

- Accuracy
- Imbalanced Data, Precision, Recall

Adjusting the Classification Threshold

- A case study
- ROC curves, and AUC

[Extra] Detailed MLE, Gradient Descent, PR curves

# Next Time

|  | Regression ($y \in \mathbb{R}$) | Classification ($y \in \{0, 1\}$) |
|---|---|---|

**1. Choose a model** ✅

Linear Regression
$$\hat{y} = f_\theta(x) = x^T\theta$$

Logistic Regression
$$\hat{P}_\theta(Y = 1|x) = \sigma(x^T\theta)$$

**2. Choose a loss function** ✅

Squared Loss or Absolute Loss

Average Cross-Entropy Loss
$$-\frac{1}{n}\sum_{i=1}^{n}\left(y_i\log(\sigma(X_i^T\theta) + (1 - y_i)\log(1 - \sigma(X_i^T\theta))\right)$$

**3. Fit the model**

Regularization
Sklearn/Gradient descent

Regularization
Sklearn/~~Gradient descent~~

**4. Evaluate model performance**

$R^2$, Residuals, etc.

Let's do it!

# Classifier Accuracy

Now that we actually have our classifier, let's try and quantify how well it performs.

The most basic evaluation metric for a classifier is **accuracy**.

$$\text{accuracy} = \frac{\#\text{ of points classified correctly}}{\#\text{ points total}}$$

```python
def accuracy(X, Y):
    return np.mean(model.predict(X) == Y)

accuracy(X, Y) # 0.8691
```

```python
model.score(X, Y) # 0.8691
```
(sklearn [documentation](#))

While widely used, the accuracy metric is **not so meaningful** when dealing with **class imbalance** in a dataset.

# Pitfalls of Accuracy: A Case Study

Suppose we're trying to build a classifier to filter spam emails.

- Each email is **spam** (1) or **ham** (0).

Let's say we have 100 emails, of which only **5** are truly **spam**, and the remaining **95** are **ham**.

Your friend ("Friend 1"):

Classify every email as **ham** (0).

$$\hat{y} = \text{classify}_{\text{friend}}(x) = 0$$

1. What is the accuracy of your friend's classifier?
2. Is accuracy a good metric of this classifier's performance?

# Pitfalls of Accuracy: A Case Study

Suppose we're trying to build a classifier to filter spam emails.

- Each email is **spam** (1) or **ham** (0).

Let's say we have 100 emails, of which only **5** are truly **spam**, and the remaining **95** are **ham**.

Your friend ("Friend 1"):

Classify every email as **ham** (0).

$$\text{accuracy}_1 = \frac{95}{100} = 0.95$$

**High** accuracy…
…but we detected **none** ⚠️ of the spam!!!

# Pitfalls of Accuracy: A Case Study

Suppose we're trying to build a classifier to filter spam emails.

- Each email is **spam** (1) or **ham** (0).

Let's say we have 100 emails, of which only **5** are truly **spam**, and the remaining **95** are **ham**.

Your friend ("Friend 1"):

Classify every email as **ham** (0).

$$\text{accuracy}_1 = \frac{95}{100} = 0.95$$

**High** accuracy…
…but we detected **none** ⚠ of the spam!!!

Your other friend ("Friend 2"):

Classify every email as **spam** (1).

$$\text{accuracy}_2 = \frac{5}{100} = 0.05$$

**Low** ⚠ accuracy…
…but we detected **all** of the spam!!!

# Pitfalls of Accuracy: Class Imbalance

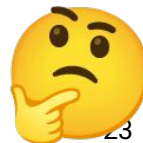Suppose we're trying to build a classifier to filter spam emails.

- Each email is **spam** (1) or **ham** (0).

Let's say we have 100 emails, of which only **5** are truly **spam**, and the remaining **95** are **ham**.

> Accuracy is not always a good metric for classification, particularly when your data have **class imbalance** (e.g., very few 1's compared to 0's).

Your friend ("Friend 1"):

Classify every email as **ham** (0).

$$\text{accuracy}_1 = \frac{95}{100} = 0.95$$

**High** accuracy…
…but we detected **none** ⚠️ of the spam!!!

Your other friend:

Classify every email as **spam** (1).

$$\text{accuracy}_2 = \frac{5}{100} = 0.05$$

**Low** ⚠️ accuracy…
…but we detected **all** of the spam!!!

# Types of Classification Successes/Errors: The Confusion Matrix

- **True positives** and **true negatives** are when we correctly classify an observation as being positive or negative, respectively.

- **False positives** are "false alarms": we predicted 1, but the true class was 0.
- **False negatives** are "failed detections": we predicted 0, but the true class was 1.

Prediction $\hat{y}$

| Actual $y$ | | **0** | **1** |
|---|---|---|---|
| | 0 | True **negative** (TN) | False **positive** (FP) |
| | 1 | False **negative** (FN) | True **positive** (TP) |

"**positive**" means a prediction of **1**.
"**negative**" means a prediction of **0**.

# Types of Classification Successes/Errors: The Confusion Matrix

- **True positives** and **true negatives** are when we correctly classify an observation as being positive or negative, respectively.

- **False positives** are "false alarms": we predicted 1, but the true class was 0.
- **False negatives** are "failed detections": we predicted 0, but the true class was 1.

Prediction $\hat{y}$

|  | 0 | 1 |
|---|---|---|
| Actual $y$   0 | True **negative** (TN) | False **positive** (FP) |
| 1 | False **negative** (FN) | True **positive** (TP) |

A confusion matrix plots these four quantities for a particular classifier and dataset.

```
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(Y_true, Y_pred)
```

# Accuracy, Precision, and Recall

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

What proportion of points did our classifier classify correctly?

| | Prediction | |
|---|---|---|
| | **0** | **1** |
| 0 | TN | FP |
| 1 | FN | TP |

Actual

# Accuracy, Precision, and Recall

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

What proportion of points did our classifier classify correctly?

| | | Prediction | |
|---|---|---|---|
| | | **0** | **1** |
| Actual | 0 | TN | FP |
| | 1 | FN | TP |

**Precision** and recall are two commonly used metrics that,
measure performance even in the presence of class imbalance.

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

Of all observations that were predicted to be 1, what proportion were actually 1?
- How **accurate** is our classifier **when it is positive**?
- Penalizes false positives.

# Accuracy, Precision, and Recall

| | | Prediction | |
|---|---|---|---|
| | | **0** | **1** |
| Actual | 0 | TN | FP |
| | 1 | FN | TP |

$$accuracy = \frac{TP + TN}{n}$$

What proportion of points did our classifier classify correctly?

Precision and **recall** are two commonly used metrics that, measure performance even in the presence of class imbalance.

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

Of all observations that were predicted to be 1, what proportion were actually 1?
- How accurate is our classifier when it is positive?
- Penalizes false positives.

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

Of all observations that were actually 1, what proportion did we predict to be 1? (Also known as sensitivity.)
- How **sensitive** is our classifier to **positives**?
- Penalizes false negatives.

# One of the Most Valuable Graphics on Wikipedia

(i.e., positive;
predicted class is 1)



relevant elements

| false negatives | true negatives |
| true positives | false positives |

retrieved elements

(*i.e., true class is 1)

How many retrieved items are relevant?

$$\text{Precision} = \frac{}{}$$

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

How many relevant items are retrieved?

$$\text{Recall} = \frac{}{}$$

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

[adapted from Wikipedia]

# Back to the Spam

Suppose we're trying to build a classifier to filter spam emails.

- Each email is **spam** (1) or **ham** (0).

Let's say we have 100 emails, of which only **5** are truly **spam**, and the remaining **95** are **ham**.

Your friend:

Classify every email as **ham** (0).

|   | **0** | **1** |
|---|-------|-------|
| 0 | TN: 95 | FP: 0 |
| 1 | FN: 5 | TP: 0 |

$$\text{accuracy}_1 = \frac{95}{100} = 0.95$$

$$\text{precision}_1 = \frac{0}{0 + 0} = \text{undefined}$$

$$\text{recall}_1 = \frac{0}{0 + 5} = 0$$

33

# Back to the Spam

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

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

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

Suppose we're trying to build a classifier to filter spam emails.

- Each email is **spam** (1) or **ham** (0).

Let's say we have 100 emails, of which only **5** are truly **spam**, and the remaining **95** are **ham**.

Your friend:

Classify every email as **ham** (0).

$$\text{accuracy}_1 = \frac{95}{100} = 0.95$$

Never positive!
$$\begin{cases} \text{precision}_1 = \frac{0}{0 + 0} = \text{undefined} \\ \\ \text{recall}_1 = \frac{0}{0 + 5} = 0 \end{cases}$$

Your other friend ("Friend 2"):

Classify every email as **spam** (1).

|   | **0** | **1** |
|---|---|---|
| 0 | TN: 0 | FP: 95 |
| 1 | FN: 0 | TP: 5 |

$$\text{accuracy}_2 = \frac{5}{100} = 0.05$$

$$\text{precision}_2 = \frac{5}{5 + 95} = 0.05$$ } Many false positives!

$$\text{recall}_2 = \frac{5}{5 + 0} = 1.0$$ } No false negatives!

# Precision vs. Recall

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

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

Precision penalizes false positives, and Recall penalizes false negatives.

We can achieve **100% recall** by making our classifier output "1", regardless of the input.

- Friend 2's "always predict spam" classifier.
- We would have no false negatives, but many false positives, and so our **precision would be low**.

This suggests that there is a **tradeoff** between precision and recall; they are often inversely related.

- Ideally, both would be near 100%, but that's unlikely to happen.

(see extra slides re: the precision-recall curve)

# Which Performance Metric?

$$accuracy = \frac{TP + TN}{n}$$

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

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

In many settings, there might be a much higher cost to missing positive cases.
For our tumor classifier:

- We really don't want to miss any malignant tumors (avoid false negatives).
- We might be fine with classifying benign tumors as malignant (OK to have false positives), since pathologists could do further studies to verify all malignant tumors.
- This context would prioritize **recall**.

How do we engineer classifiers to meet the performance goals of our problem?

36

PCA clustering Principal Component Analysis (PCA) plot of 20 populations from 1000 Genomes Project, built using 2 first principal components. The following populations were not used to build the

# Interlude

PCA is commonly used in biomedical contexts, which have many named variables!

1. To cluster data (Paper 1, Paper 2)



2. To identify correlated variables (interpret rows of $V^T$ as linear coefficients) (Paper 3). Uses biplots.



**Break (2 min)**

37

# Adjusting the Classification Threshold

Lecture 22, Data 100 Spring 2022

Logistic Regression Model, continued

- sklearn demo
- Maximum Likelihood Estimation: high-level (live), detailed (recorded)

Linear separability and Regularization

Performance Metrics

- Accuracy
- Imbalanced Data, Precision, Recall

**Adjusting the Classification Threshold**

- A case study
- ROC curves, and AUC

[Extra] Detailed MLE, Gradient Descent, PR curves

**Feature Engineering**:

What are the features x that generate great probabilities for prediction?

Parameter $\theta = (\theta_0, \theta_1, \ldots, \theta_p)$ :

$$x \longrightarrow \boxed{\hat{P}_\theta(Y = 1|x) = \sigma(x^T\theta)} \longrightarrow$$ Probability of response = 1

observation vector

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

classify

prediction    **categorical** response

**Classification**:

What is the best
**classification threshold T** to
choose that best fits our
problem context?

Parameter $\theta = (\theta_0, \theta_1, \ldots, \theta_p)$ :

$$x \longrightarrow \boxed{\hat{P}_\theta(Y = 1|x) = \sigma(x^T\theta)} \longrightarrow$$

observation
vector

Probability of
response = 1

$$\hat{y} = \text{classify}(x) = \begin{cases} 1 & \hat{P}_\theta(Y = 1|x) \geq T \\ 0 & \text{otherwise} \end{cases}$$

sklearn's
model.predict()
uses fixed 0.5

classify

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

prediction          **categorical**
                    response

40

# Classification Threshold

$$\hat{y} = \text{classify}(x) = \begin{cases} 1 & \hat{P}_\theta(Y = 1|x) \geq T \\ 0 & \text{otherwise} \end{cases}$$

The default threshold in sklearn is T = 0.5.



True classes and model fit

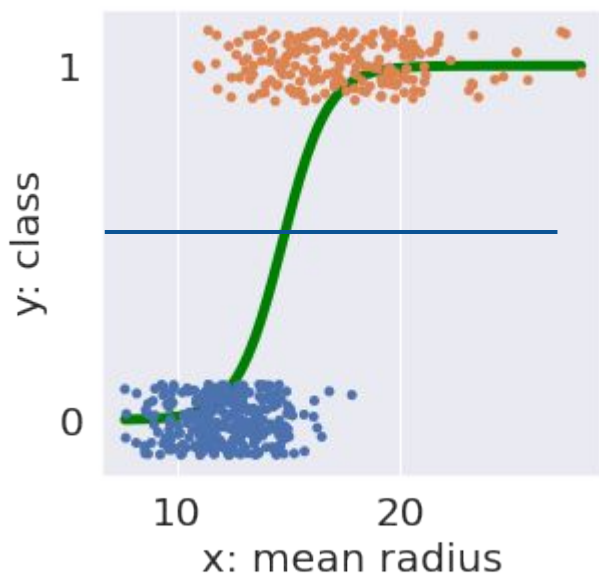| | X | Y | P(Y = 1 | x) |
|---|---|---|---|
| 0 | 25.220 | 1 | 0.999965 |
| 1 | 13.480 | 1 | 0.226448 |
| 2 | 11.290 | 0 | 0.033174 |
| 3 | 12.860 | 0 | 0.137598 |
| 4 | 19.690 | 1 | 0.992236 |
| ... | ... | ... | ... |
| 507 | 8.888 | 0 | 0.003257 |
| 508 | 11.640 | 0 | 0.046105 |
| 509 | 14.290 | 0 | 0.392796 |
| 510 | 13.980 | 1 | 0.323216 |
| 511 | 12.180 | 0 | 0.075786 |

Predicted classes if T = 0.5

41

# Classification Threshold

$$\hat{y} = \mathrm{classify}(x) = \begin{cases} 1 & \hat{P}_\theta(Y = 1|x) \geq T \\ 0 & \mathrm{otherwise} \end{cases}$$
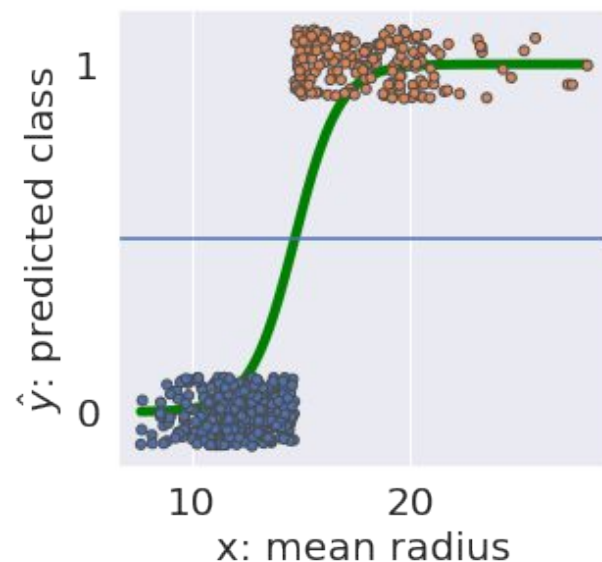
As we increase the threshold T, we "raise the standard" of how confident our classifier needs to be to predict 1 (i.e., "positive").

T = 0.25

T = 0.50

T = 0.75



These x will all predict 1

Fewer positives

## Choosing an Accuracy Threshold

The choice of threshold T impacts our classification performance.

- High T:     Most predictions are 0. Lots of false negatives.
- Low T:      Most predictions are 1. Lots of false positives.

Do we get max accuracy when T ≈ 0.5? Not always the case…

**Demo**

See notebook for code snippets.

Train Accuracy vs. Threshold



T ≈ 0.57

Best T ≈ 0.57 likely due to class imbalance. There are fewer malgnant tumors and so we want to be more confident before classifying a tumor as malignant.

```
malignant
0           317
1           195
dtype: int64
```

The threshold should typically be tuned using cross validation.

For a threshold T:

$$\text{cross\_val\_acc} = (1/k) \left[ \boxed{\begin{array}{c} \text{Model fit to} \\ \text{train set 1,} \\ \text{Acc on val set} \\ 1 \end{array}} + ... + \boxed{\begin{array}{c} \text{Model fit to} \\ \text{train set k,} \\ \text{Acc on val set} \\ k \end{array}} \right]$$

Cross-Validated Accuracy vs. Threshold



T ≈ 0.56

**Demo**

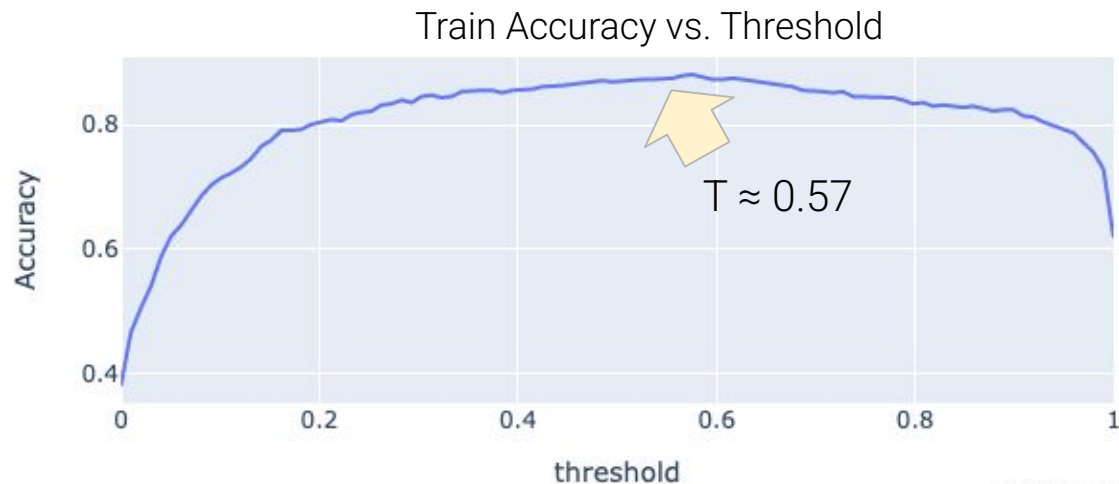See notebook for code snippets.
documentation

## Choosing a Threshold According to Other Metrics?

The choice of threshold T impacts our classification performance.

- High T: Most predictions are 0. Lots of false negatives.
- Low T: Most predictions are 1. Lots of false positives.

---

Could we choose a threshold T based on metrics that measure false positives/false negatives?

Yes! Two options:

- Precision-Recall Curve (PR Curve). Covered in extra slides.
- "Receiver Operating Characteristic" Curve (**ROC Curve**).

Each of these visualizations have an associated performance metric: **AUC (Area Under Curve)**.

**Demo**

**Demo**

## Two More Metrics

$$\mathrm{TPR} = \frac{TP}{TP + FN}$$

**True Positive Rate** (TPR):
"What proportion of spam did I mark correctly?
Same thing as **recall**. In statistics, sensitivity.

$$\mathrm{FPR} = \frac{FP}{FP + TN}$$

**False Positive Rate** (FPR):
"What proportion of regular email did I mark as spam?
In statistics, also called specificity.

| | | Prediction | |
|---|---|---|---|
| | | 0 | 1 |
| 0 | | TN | FP |
| 1 | | FN | TP |

The ROC curve plots TPR vs FPR for different classification thresholds.

relevant elements

false negatives     true negatives

true positives     false positives

retrieved elements

retrieved elements

How many retrieved items are relevant?

$$\text{Precision} = \frac{\phantom{TP}}{\phantom{TP+FP}}$$

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

Want close to 1.0

How many relevant items are retrieved?

$$\text{Recall} = \frac{\phantom{TP}}{\phantom{TP+FN}}$$

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

Want close to 1.0

How many irrelevant items are retrieved?

$$\textbf{FPR} = \frac{\phantom{FP}}{\phantom{FP+TN}}$$

$$\text{FPR} = \frac{FP}{FP + TN}$$

Want close to 0.0

The ROC curve plots TPR vs FPR for different classification thresholds.

47

The choice of threshold T impacts our classification performance.

- High T:    Most predictions are 0. Lots of false negatives.
- Low T:     Most predictions are 1. Lots of false positives.

$$\text{TPR} = \frac{TP}{TP + FN} \qquad \qquad \text{FPR} = \frac{FP}{FP + TN}$$

**Demo**



As we increase T, **both TPR** and FPR decrease.

- A decreased **TPR** is bad (detecting fewer positives).
- A decreased **FPR** is good (fewer false positives).

48

## The ROC Curve

$$\mathrm{TPR} = \frac{TP}{TP + FN}$$

$$\mathrm{FPR} = \frac{FP}{FP + TN}$$

The ROC Curve plots this tradeoff.

- ROC stands for "Receiver Operating Characteristic." [Wikipedia]
- We want high TPR, low FPR.



**Demo**

1. Which part of this curve corresponds to T = 0.9?
2. Which part of this curve corresponds to T = 0.1?

$$\mathrm{TPR} = \frac{TP}{TP + FN}$$

$$\mathrm{FPR} = \frac{FP}{FP + TN}$$

The ROC Curve plots this tradeoff.

- ROC stands for "Receiver Operating Characteristic." [Wikipedia]
- We want high TPR, low FPR.

**Demo**

T = 0.1

T = 0.6

T = 0.9

## The Perfect Classifier

The "perfect" classifier is the one that has a
TPR of 1, and FPR of 0.

- We want our logistic regression model to match that as well as possible.
- We want our ROC curve to be as close to the "top left" of this graph as possible.

**Demo**



Perfect Predictor

# Performance Metric: Area Under Curve (AUC)

$$\mathrm{TPR} = \frac{TP}{TP + FN}$$

$$\mathrm{FPR} = \frac{FP}{FP + TN}$$

The "perfect" classifier is the one that has a
TPR of 1, and FPR of 0.

- We want our model to match
  that as well as possible.
- We want our ROC curve to be
  as close to the "top left" of this
  graph as possible.



**Perfect Predictor**

We can compute the **area under curve (AUC)** of our model.

- Different AUCs for both ROC curves and PR curves, but
  ROC is more common.
- Best possible AUC = 1. Terrible AUC = 0.5.
  - Random predictors have an AUC of around 0.5.
    Why?
- Your model's AUC: somewhere between 0.5 and 1.

52

# [Extra] What is the "worst" AUC and why is it 0.5?

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

$$\text{FPR} = \frac{FP}{FP + TN}$$

Best possible AUC = 1. Terrible AUC = 0.5.

- Random predictors have an AUC of around 0.5. Why?

A **random predictor** randomly predicts $P(Y = 1 \mid x)$ to be uniformly between 0 and 1.

On average, if your dataset is size $n_1 + n_0$ (with $n_1$ true class 1's and $n_0$ true class 0's):

If T = 0.5:

|   | 0 | 1 |
|---|---|---|
| 0 | TN = 0.5 $n_0$ | FP = 0.5 $n_0$ |
| 1 | FN = 0.5 $n_1$ | TP = 0.5 $n_1$ |

FPR = 0.5 $n_0$/((0.5 + 0.5)$n_0$) = 0.5

TPR = 0.5 $n_1$/((0.5 + 0.5)$n_1$) = 0.5

Point on ROC curve is (0.5, 0.5).

If T = 0.8:

|   | 0 | 1 |
|---|---|---|
| 0 | TN = 0.8 $n_0$ | FP = 0.2 $n_0$ |
| 1 | FN = 0.8 $n_1$ | TP = 0.2 $n_1$ |

FPR = 0.2 $n_0$/((0.2 + 0.8)$n_0$) = 0.2

TPR = 0.2 $n_1$/((0.2 + 0.8)$n_1$) = 0.2

Point on ROC curve is (0.2, 0.2).

If T = 0.3:

|   | 0 | 1 |
|---|---|---|
| 0 | TN = 0.3 $n_0$ | FP = 0.7 $n_0$ |
| 1 | FN = 0.3 $n_1$ | TP = 0.7 $n_1$ |

FPR = 0.7 $n_0$/((0.7 + 0.3)$n_0$) = 0.7

TPR = 0.7 $n_1$/((0.7 + 0.3)$n_1$) = 0.7

Point on ROC curve is (0.7, 0.7).

# [Extra] What is the "worst" AUC and why is it 0.5?

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

$$\text{FPR} = \frac{FP}{FP + TN}$$

Best possible AUC = 1. Terrible AUC = 0.5.

- Random predictors have an AUC of around 0.5. Why?

A **random predictor** randomly predicts P(Y = 1 | x) to be uniformly between 0 and 1.

**This slide was added post-lecture to clarify closing comments.**

**Perfect Predictor. AUC = 1.0**



**Random Predictor. AUC = 0.5.**

Area Under Curve (AUC) of random predictor is the area of this triangle.

# Common techniques for evaluating classifiers

Numerical assessments:

- **Accuracy, precision, recall/TPR, FPR.**
- Area under curve (AUC), for ROC curves.

Visualizations:

- **Confusion matrices.**
- Precision/recall curves.
- ROC curves.



Terminology and derivations from a confusion matrix:

**condition positive (P)**
the number of real positive cases in the data
**condition negative (N)**
the number of real negative cases in the data

**true positive (TP)**
eqv. with hit
**true negative (TN)**
eqv. with correct rejection
**false positive (FP)**
eqv. with false alarm, Type I error
**false negative (FN)**
eqv. with miss, Type II error

**sensitivity, recall, hit rate, or true positive rate (TPR)**
$$TPR = \frac{TP}{P} = \frac{TP}{TP + FN} = 1 - FNR$$
**specificity, selectivity or true negative rate (TNR)**
$$TNR = \frac{TN}{N} = \frac{TN}{TN + FP} = 1 - FPR$$
**precision or positive predictive value (PPV)**
$$PPV = \frac{TP}{TP + FP} = 1 - FDR$$
**negative predictive value (NPV)**
$$NPV = \frac{TN}{TN + FN} = 1 - FOR$$
**miss rate or false negative rate (FNR)**
$$FNR = \frac{FN}{P} = \frac{FN}{FN + TP} = 1 - TPR$$
**fall-out or false positive rate (FPR)**
$$FPR = \frac{FP}{N} = \frac{FP}{FP + TN} = 1 - TNR$$
**false discovery rate (FDR)**
$$FDR = \frac{FP}{FP + TP} = 1 - PPV$$
**false omission rate (FOR)**
$$FOR = \frac{FN}{FN + TN} = 1 - NPV$$
**Threat score (TS) or Critical Success Index (CSI)**
$$TS = \frac{TP}{TP + FN + FP}$$

**accuracy (ACC)**
$$ACC = \frac{TP + TN}{P + N} = \frac{TP + TN}{TP + TN + FP + FN}$$
**F1 score**
is the harmonic mean of precision and sensitivity
$$F_1 = 2 \cdot \frac{PPV \cdot TPR}{PPV + TPR} = \frac{2TP}{2TP + FP + FN}$$
**Matthews correlation coefficient (MCC)**
$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$
**Informedness or Bookmaker Informedness (BM)**
$$BM = TPR + TNR - 1$$
**Markedness (MK)**
$$MK = PPV + NPV - 1$$

We're only scratching the surface here.

55

# Extra Slides

Lecture 22, Data 100 Spring 2022

Logistic Regression Model, continued

- sklearn demo
- Maximum Likelihood Estimation: high-level (live), detailed (recorded)

Linear separability and Regularization

Performance Metrics

- Accuracy
- Imbalanced Data, Precision, Recall

Adjusting the Classification Threshold

- A case study
- ROC curves, and AUC

**[Extra] Detailed MLE, Gradient Descent, PR curves**

**Video**: link
Out of scope, but useful for understanding where cross-entropy loss comes from.

# [Extra] Detailed Maximum Likelihood Estimation

Lecture 22, Data 100 Spring 2022

Logistic Regression Model, continued

- sklearn demo
- Maximum Likelihood Estimation: high-level (live), detailed (recorded)

Linear separability and Regularization

Performance Metrics

- Accuracy
- Imbalanced Data, Precision, Recall

Adjusting the Classification Threshold

- A case study
- ROC curves, and AUC

**[Extra] Detailed MLE**, Gradient Descent, PR curves

# The Modeling Process

|  | Regression ($y \in \mathbb{R}$) | Classification ($y \in \{0, 1\}$) |
|---|---|---|
| **1. Choose a model** ✅ | Linear Regression $\hat{y} = f_\theta(x) = x^T\theta$ | Logistic Regression $\hat{P}_\theta(Y = 1|x) = \sigma(x^T\theta)$ |
| **2. Choose a loss function** | Squared Loss or Absolute Loss | Average Cross-Entropy Loss $-\frac{1}{n}\sum_{i=1}^{n}\left(y_i\log(\sigma(X_i^T\theta) + (1 - y_i)\log(1 - \sigma(X_i^T\theta))\right)$ |
| **3. Fit the model** | Regularization Sklearn/Gradient descent | Regularization Sklearn/Gradient descent |
| **4. Evaluate model performance** | $R^2$, Residuals, etc. | ?? (next time) |

Wherefore use cross-entropy?

Shakespeare [Wikipedia]

# Why Use Cross-Entropy Loss?

This section will not be directly tested, but you will understand why we minimize cross-entropy loss for logistic regression.

Two common explanations:
- [Information Theory] KL Divergence ([textbook](#))
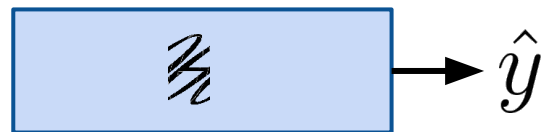- [Probability] Maximum Likelihood Estimation (this lecture)

# Recall the Coin Demo (No-Input Classification)

For training data:        {0, 0, 1, 1, 1, 1, 0, 0, 0, 0}

0.4 is the most "intuitive" θ for two reasons:

1. Frequency of heads in our data
2. Maximizes the **likelihood** of our data:  (**proportional to** the probability of our data)

$$\hat{\theta} = \operatorname*{argmax}_{\theta} \left( \theta^4 (1 - \theta)^6 \right)$$



Parameter θ:
Probability that
IID flip == 1 (Heads)

Prediction:
1 or 0

How can we generalize this notion of likelihood to **any** random binary sample?

$$\{y_1, y_2, \ldots, y_n\} \implies \hat{\theta} = \operatorname*{argmax}_{\theta} (???)$$

data (1's and 0's)        likelihood

# A Compact Representation of the Bernoulli Probability Distribution

How can we generalize this notion of likelihood to **any** random binary sample?

$$\{y_1, y_2, \ldots, y_n\} \implies \hat{\theta} = \underset{\theta}{\mathrm{argmax}} \, (???)$$

data (1's and 0's)    likelihood

Let $Y$ be Bernoulli($p$). The probability distribution can be written compactly:

$$P(Y = y) = p^y (1 - p)^{1-y}$$

For P(Y = **1**), only this term stays

For P(Y = **0**), only this term stays

(long, non-compact form):

$$P(Y = y) = \begin{cases} p & \text{if } y = 1 \\ 1 - p & \text{if } y = 0 \end{cases}$$

# Generalized Likelihood of Binary Data

How can we generalize this notion of likelihood to **any** random binary sample?

$$\{y_1, y_2, \ldots, y_n\} \implies \hat{\theta} = \underset{\theta}{\mathrm{argmax}}\,(???)$$

data (1's and 0's)　　　　　　　　likelihood

Let $Y$ be Bernoulli($p$). The probability distribution can be written compactly:

$$P(Y = y) = p^y (1 - p)^{1-y}$$

For P(Y = **1**), only this term stays　　For P(Y = **0**), only this term stays

If binary data are **IID with same** probability p, then the likelihood of the data is:

$$\prod_{i=1}^{n} p^{y_i} (1 - p)^{(1 - y_i)}$$

Ex: $\{0, 0, 1, 1, 1, 1, 0, 0, 0, 0\} \rightarrow p^4 (1 - p)^6$

likelihood vs. probability

$$\binom{10}{4} p^4 (1 - p)^6$$

# Generalized Likelihood of Binary Data

How can we generalize this notion of likelihood to any random binary sample?

$$\{y_1, y_2, \ldots, y_n\} \implies \hat{\theta} = \underset{\theta}{\operatorname{argmax}} \, (???)$$

data (1's and 0's)                                                 likelihood

Let $Y$ be Bernoulli($p$). The probability distribution can be written compactly:

$$P(Y = y) = p^y (1-p)^{1-y}$$

For P(Y = **1**), only this term stays

For P(Y = **0**), only this term stays

If binary data are **IID with same** probability p, then the likelihood of the data is:

$$\prod_{i=1}^{n} p^{y_i} (1-p)^{(1-y_i)}$$

If binary data are independent with **different** probability $p_i$, then the likelihood of the data is:

(spoiler: for logistic regression, $p_i = \sigma(X_i^T \theta)$)

$$\prod_{i=1}^{n} p_i^{y_i} (1-p_i)^{(1-y_i)}$$

# Maximum Likelihood Estimation (MLE)

Our **maximum likelihood estimation** problem:

- For i = 1, 2, …, n, let $Y_i$ be independent Bernoulli($p_i$). Observe data $\{y_1, y_2, \ldots, y_n\}$.
- We'd like to estimate $p_1, p_2, \ldots, p_n$.

Find $\hat{p}_1, \hat{p}_2, \ldots, \hat{p}_n$ that **maximize** $\displaystyle\prod_{i=1}^{n} p_i^{y_i} (1 - p_i)^{(1 - y_i)}$

# Maximum Likelihood Estimation (MLE)

Our **maximum likelihood estimation** problem:

- For i = 1, 2, …, n, let $Y_i$ be independent Bernoulli($p_i$). Observe data $\{y_1, y_2, \ldots, y_n\}$.
- We'd like to estimate $p_1, p_2, \ldots, p_n$.

Find $\hat{p}_1, \hat{p}_2, \ldots, \hat{p}_n$ that **maximize** $\displaystyle\prod_{i=1}^{n} p_i^{y_i} (1 - p_i)^{(1-y_i)}$

Equivalent, simplifying optimization problems:

$$\underset{p_1, p_2, \ldots, p_n}{\textbf{maximize}} \quad \log\left(\prod_{i=1}^{n} p_i^{y_i} (1 - p_i)^{(1-y_i)}\right) \qquad \text{(log is an increasing function. If a > b, then log(a) > log(b).)}$$
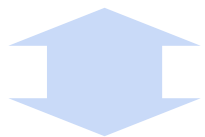
$$= \sum_{i=1}^{n} \log\left(p_i^{y_i} (1 - p_i)^{(1-y_i)}\right) = \sum_{i=1}^{n} \left(y_i \log(p_i) + (1 - y_i) \log(1 - p_i)\right)$$

# Maximum Likelihood Estimation (MLE)

Our **maximum likelihood estimation** problem:

- For i = 1, 2, …, n, let $Y_i$ be independent Bernoulli($p_i$). Observe data $\{y_1, y_2, \ldots, y_n\}$.
- We'd like to estimate $p_1, p_2, \ldots, p_n$.

Find $\hat{p}_1, \hat{p}_2, \ldots, \hat{p}_n$ that **maximize** $\displaystyle\prod_{i=1}^{n} p_i^{y_i}(1-p_i)^{(1-y_i)}$

Equivalent, simplifying optimization problems:

$$\underset{p_1, p_2, \ldots, p_n}{\textbf{maximize}} \quad \log\left(\prod_{i=1}^{n} p_i^{y_i}(1-p_i)^{(1-y_i)}\right) \qquad \text{(log is an increasing function. If a > b, then log(a) > log(b).)}$$

$$= \sum_{i=1}^{n} \log\left(p_i^{y_i}(1-p_i)^{(1-y_i)}\right) = \sum_{i=1}^{n}\left(y_i \log(p_i) + (1-y_i)\log(1-p_i)\right)$$

$$\underset{p_1, p_2, \ldots, p_n}{\textbf{minimize}} \quad -\frac{1}{n}\sum_{i=1}^{n}\left(y_i \log(p_i) + (1-y_i)\log(1-p_i)\right)$$

Argmax property:
x that maximizes f(x) will
minimize -f(x)

66

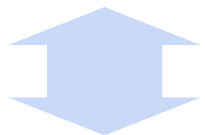# Maximizing Likelihood == Minimizing Average Cross-Entropy

**argmax** $p_1, p_2, \ldots, p_n$   $\displaystyle\prod_{i=1}^{n} p_i^{y_i}(1-p_i)^{(1-y_i)}$

Log is increasing;
max/min properties

**argmin** $p_1, p_2, \ldots, p_n$   $\displaystyle -\frac{1}{n}\sum_{i=1}^{n}\left(y_i \log(p_i) + (1-y_i)\log(1-p_i)\right)$   **Average Cross-Entropy Loss!!**

For logistic regression,
let $p_i = \sigma(X_i^T \theta)$

**argmin** $\theta$   $\displaystyle -\frac{1}{n}\sum_{i=1}^{n}\left(y_i \log(\sigma(X_i^T\theta)) + (1-y_i)\log(1-\sigma(X_i^T\theta))\right)$   **Average Cross-Entropy Loss for Logistic Regression!!**

**Minimizing cross-entropy loss** is equivalent to **maximizing the likelihood of the training data**.

Assumption: all data are independent Bernoulli random variables.

$$\underset{\theta}{\textbf{argmin}} \quad -\frac{1}{n}\sum_{i=1}^{n}\left(y_i\log(\sigma(X_i^T\theta) + (1-y_i)\log(1-\sigma(X_i^T\theta)))\right)$$

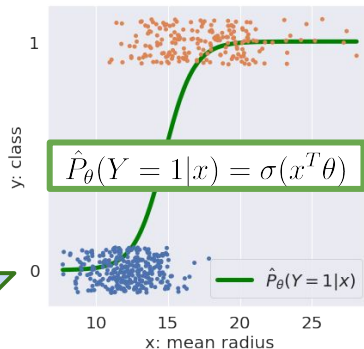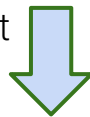For logistic regression, let $p_i = \sigma(X_i^T\theta)$

$$\underset{p_1, p_2, \ldots, p_n}{\textbf{argmax}} \quad \prod_{i=1}^{n} p_i^{y_i}(1-p_i)^{(1-y_i)}$$

Prob. that i-th response is yi

**Minimizing cross-entropy loss** is equivalent to **maximizing the likelihood of the training data**.

$$\underset{\theta}{\textbf{argmin}} \quad -\frac{1}{n}\sum_{i=1}^{n}\left(y_i \log(\sigma(X_i^T\theta) + (1-y_i)\log(1-\sigma(X_i^T\theta)))\right)$$

For logistic regression, let $p_i = \sigma(X_i^T\theta)$

Assumption: all data are independent Bernoulli random variables

$$\underset{p_1,p_2,\ldots,p_n}{\textbf{argmax}} \quad \prod_{i=1}^{n} \underbrace{p_i^{y_i}(1-p_i)^{(1-y_i)}}_{\text{Prob. that i-th response is yi}}$$

**Main takeaway**: The optimal theta that minimizes mean cross-entropy loss "pushes" all probabilities in the direction of the true class.



Want $\sigma(x^T\theta) \to 1$

$\hat{P}_\theta(Y=1|x) = \sigma(x^T\theta)$

Want $\sigma(x^T\theta) \to 0$

**Minimizing cross-entropy loss** is equivalent to **maximizing the likelihood of the training data**.

Assumption: all data are independent Bernoulli random variables

It turns out that many of the model + loss combinations we've seen can be motivated using MLE.

- OLS, Ridge Regression, etc.
- You will study MLE further in probability and ML classes. But now you know it exists.
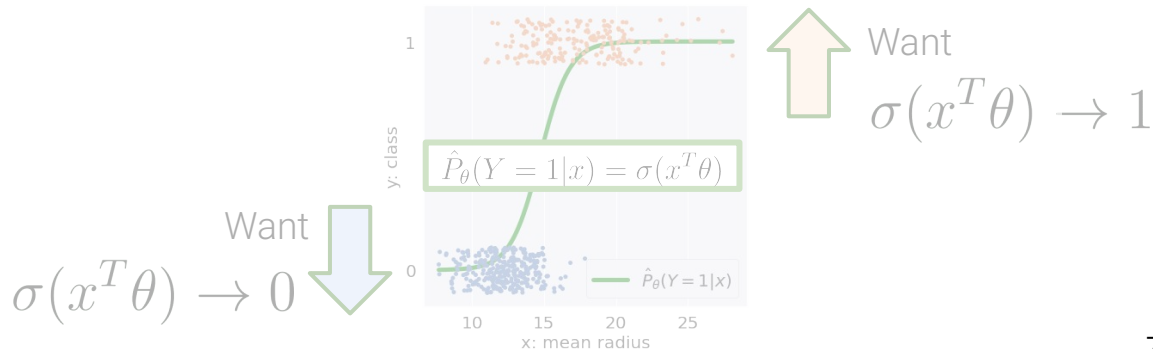
$$\underset{\theta}{\textbf{argmin}} \quad -\frac{1}{n}\sum_{i=1}^{n}\left(y_i\log(\sigma(X_i^T\theta) + (1-y_i)\log(1-\sigma(X_i^T\theta)))\right)$$

For logistic regression, let $p_i = \sigma(X_i^T\theta)$

$$\underset{p_1, p_2, \ldots, p_n}{\textbf{argmax}} \quad \prod_{i=1}^{n}\underbrace{p_i^{y_i}(1-p_i)^{(1-y_i)}}$$

Prob. that i-th response is yi



Want

$\sigma(x^T\theta) \to 1$

$\hat{P}_\theta(Y = 1|x) = \sigma(x^T\theta)$

Want

$\sigma(x^T\theta) \to 0$

# We Did it!

| | Regression ($y \in \mathbb{R}$) | Classification ($y \in \{0, 1\}$) |
|---|---|---|

**1. Choose a model** ✅

Linear Regression
$$\hat{y} = f_\theta(x) = x^T \theta$$

Logistic Regression
$$\hat{P}_\theta(Y = 1 | x) = \sigma(x^T \theta)$$

**2. Choose a loss function** ✅

Squared Loss or Absolute Loss

Average Cross-Entropy Loss
$$-\frac{1}{n} \sum_{i=1}^{n} \left( y_i \log(\sigma(X_i^T \theta) + (1 - y_i) \log(1 - \sigma(X_i^T \theta)) \right)$$

**3. Fit the model**

Regularization
Sklearn/Gradient descent

Regularization
Gradient descent

That which we call a rose would by any other name smell as sweet.

**4. Evaluate model performance**

$R^2$, Residuals, etc.

??
(next time)

Shakespeare
[]

Reference slides. Out of scope.

# [Extra] Gradient Descent for Logistic Regression

Lecture 22, Data 100 Spring 2022

Logistic Regression Model, continued

- sklearn demo
- Maximum Likelihood Estimation: high-level (live), detailed (recorded)

Linear separability and Regularization

Performance Metrics

- Accuracy
- Imbalanced Data, Precision, Recall

Adjusting the Classification Threshold

- A case study
- ROC curves, and AUC

**[Extra]** Detailed MLE, **Gradient Descent**, PR curves

|  | Regression ($y \in \mathbb{R}$) | Classification ($y \in \{0, 1\}$) |
|---|---|---|
| ✅ **1. Choose a model** | Linear Regression $$\hat{y} = f_\theta(x) = x^T\theta$$ | Logistic Regression $$\hat{P}_\theta(Y = 1|x) = \sigma(x^T\theta)$$ |
| ✅ **2. Choose a loss function** | Squared Loss or Absolute Loss | Average Cross-Entropy Loss $$-\frac{1}{n}\sum_{i=1}^{n}\left(y_i\log(\sigma(X_i^T\theta) + (1 - y_i)\log(1 - \sigma(X_i^T\theta))\right)$$ |
| **3. Fit the model** | Regularization Sklearn/Gradient descent | Regularization Sklearn/**Gradient descent** |
| ✅ **4. Evaluate model performance** | $R^2$, Residuals, etc. | Accuracy, Precision, Recall, ROC Curves |

# A Simplification

$$y_i \log(p_i) + (1 - y_i) \log(1 - p_i)$$

$$= y_i \log\left(\frac{p_i}{1 - p_i}\right) + \log(1 - p_i)$$

$$= y_i \phi(x_i)^T \theta + \log(\sigma(-\phi(x_i)^T \theta))$$

**The calculation on the left uses the following:**

$$t_i = \phi(x_i)^T \theta$$

$$p_i = \sigma(t_i)$$

$$t_i = \log\left(\frac{p_i}{1 - p_i}\right)$$

$$1 - \sigma(t_i) = \sigma(-t_i)$$

**Final form:** The best $\hat{\theta}$ is

$$\arg\min_{\theta} -\frac{1}{n} \sum_{i=1}^{n} \left(y_i \phi(x_i)^T \theta + \log\left(\sigma\left(-\phi(x_i)^T \theta\right)\right)\right)$$

➤ Want to minimize

$$\mathbf{L}(\theta) = -\frac{1}{n} \sum_{i=1}^{n} (y_i \phi(x_i)^T \theta + \log(\sigma(-\phi(x_i)^T \theta)))$$

➤ Take Derivative:

$$\nabla_\theta \mathbf{L}(\theta) = -\frac{1}{n} \sum_{i=1}^{n} \nabla_\theta y_i \phi(x_i)^T \theta + \nabla_\theta \log\left(\sigma\left(-\phi(x_i)^T \theta\right)\right)$$

$$= -\frac{1}{n} \sum_{i=1}^{n} y_i \phi(x_i) + \nabla_\theta \log\left(\sigma\left(-\phi(x_i)^T \theta\right)\right)$$

$$\nabla_\theta \mathbf{L}(\theta) = -\frac{1}{n} \sum_{i=1}^{n} y_i \phi(x_i) + \frac{1}{\sigma\left(-\phi(x_i)^T \theta\right)} \boxed{\nabla_\theta \sigma\left(-\phi(x_i)^T \theta\right)}$$

**Derivative**
$$\frac{d}{dt}\sigma(t) = \sigma(t)\sigma(-t)$$

$$= -\frac{1}{n} \sum_{i=1}^{n} y_i \phi(x_i) + \frac{\sigma\left(-\phi(x_i)^T \theta\right)}{\sigma\left(-\phi(x_i)^T \theta\right)} \sigma\left(\phi(x_i)^T \theta\right) \nabla_\theta(-\phi(x_i)^T \theta)$$

$$= -\frac{1}{n} \sum_{i=1}^{n} \left(y_i - \sigma\left(\phi(x_i)^T \theta\right)\right) \phi(x_i)$$

➤ Set derivative = 0 and solve for $\hat{\theta}$
  ➤ No general analytic solution
  ➤ Solved using numeric methods

**Gradient Descent**

$\theta^{(0)} \leftarrow$ initial vector (random, zeros …)

For $\tau$ from 0 to convergence:

$$\theta^{(\tau+1)} \leftarrow \theta^{(\tau)} - \rho(\tau)\left(\frac{1}{n}\sum_{i=1}^{n}\nabla_{\theta}\mathbf{L}_i(\theta)\Big|_{\theta=\theta^{(\tau)}}\right)$$

**Stochastic Gradient Descent**

$\theta^{(0)} \leftarrow$ initial vector (random, zeros …)

For $\tau$ from 0 to convergence:

$\mathcal{B} \sim$ Random subset of indices

$$\theta^{(\tau+1)} \leftarrow \theta^{(\tau)} - \rho(\tau)\left(\frac{1}{|\mathcal{B}|}\sum_{i\in\mathcal{B}}\nabla_{\theta}\mathbf{L}_i(\theta)\Big|_{\theta=\theta^{(\tau)}}\right)$$

Very Similar Algorithms

Assuming Decomposable Loss Functions

76

54

Reference slides. Out of scope.

# [Extra] Precision-Recall Curves

Lecture 22, Data 100 Spring 2022

Logistic Regression Model, continued

- sklearn demo
- Maximum Likelihood Estimation: high-level (live), detailed (recorded)

Linear separability and Regularization

Performance Metrics

- Accuracy
- Imbalanced Data, Precision, Recall

Adjusting the Classification Threshold

- A case study
- ROC curves, and AUC
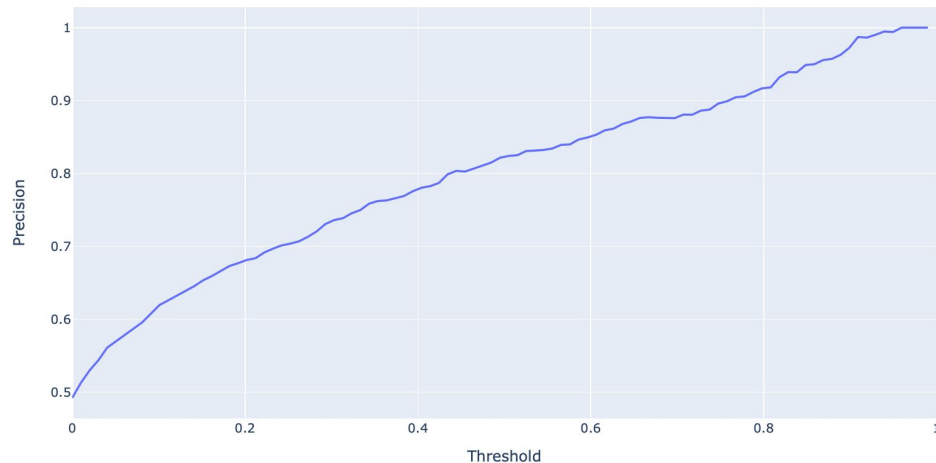
**[Extra]** Detailed MLE, Gradient Descent, **PR curves**

# Precision vs. threshold

As we increase our threshold, we have fewer and fewer false positives.

- Thus, precision tends to increase.

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

$$= \frac{\text{True Positives}}{\text{Predicted True}}$$

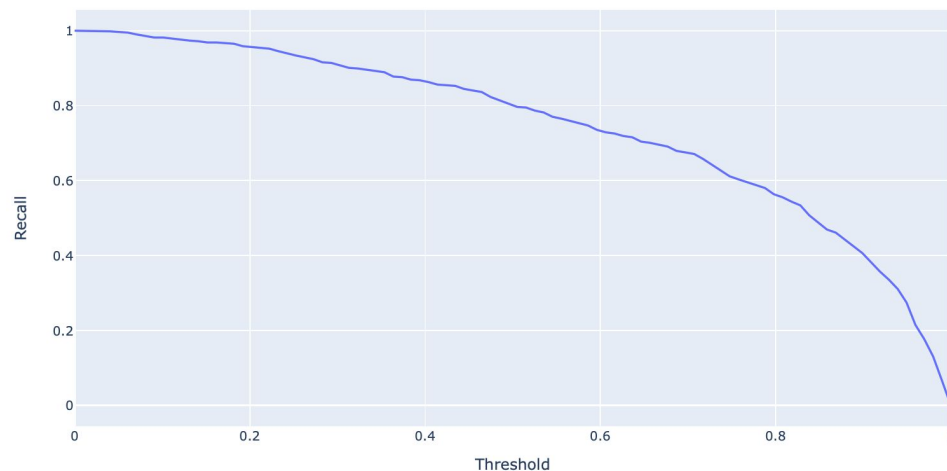It is *possible* for precision to decrease slightly with an increased threshold. Why?

# Recall vs. threshold

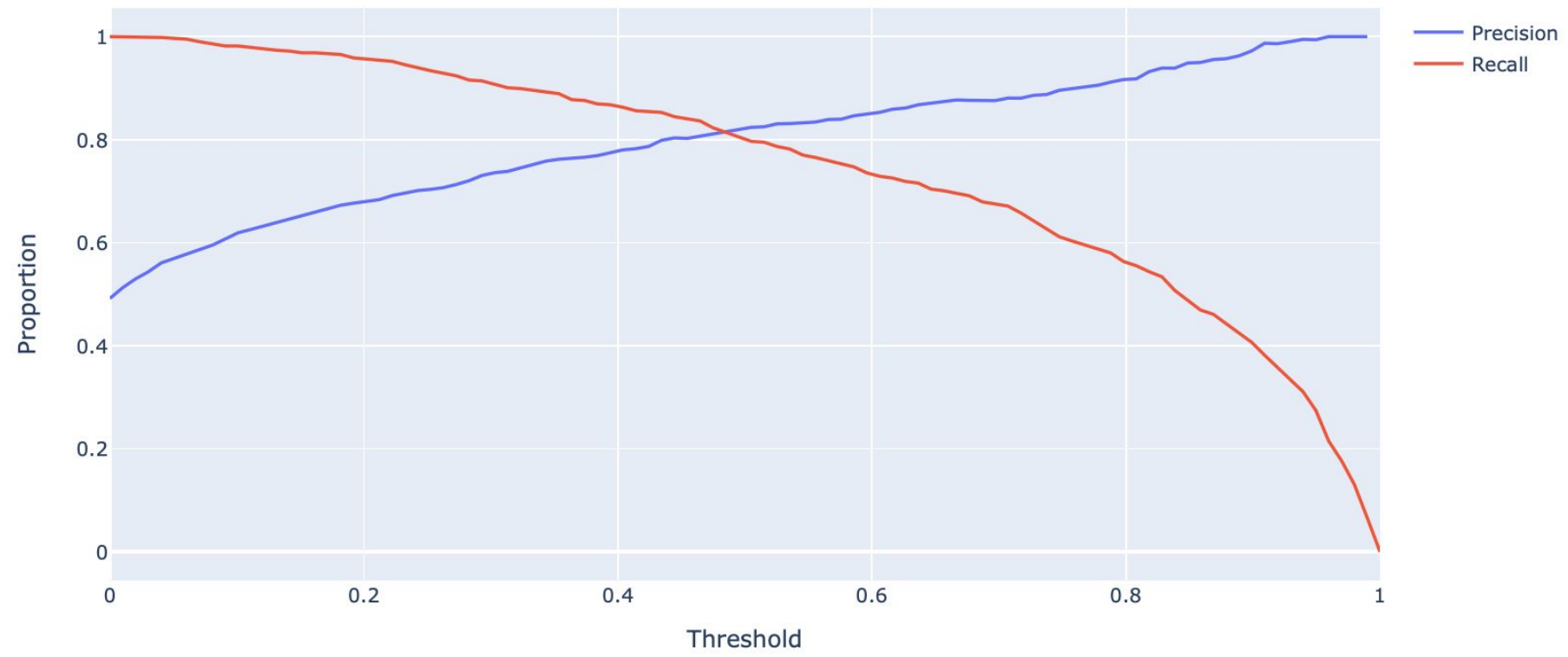As we increase our threshold, we have more and more false negatives.

- Thus, recall tends to decrease.

$$\textbf{Recall} = \frac{\textbf{True Positives}}{\textbf{True Positives} + \textbf{False Negatives}}$$

$$= \frac{\textbf{True Positives}}{\textbf{Actually True}}$$

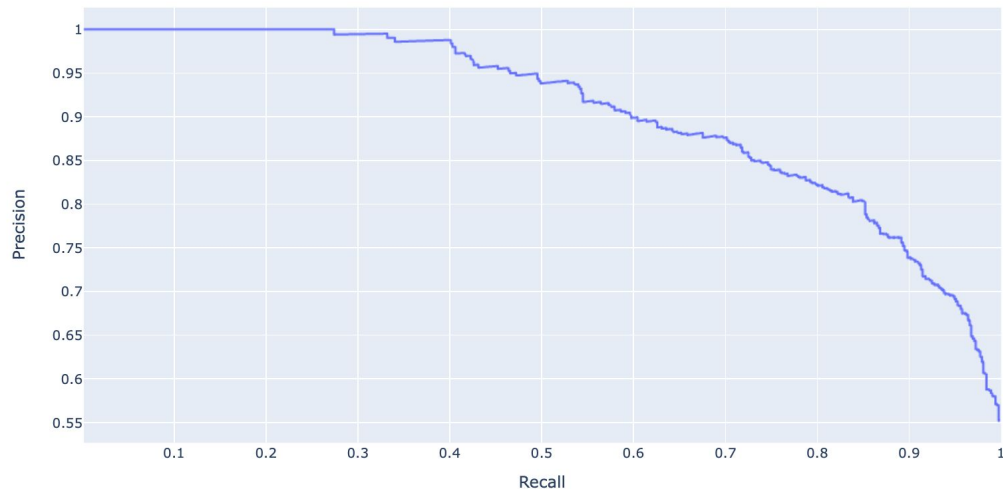Recall strictly decreases as we increase our threshold. Why?

# Precision and Recall vs. Threshold

# Precision-recall curves

We can also plot precision vs. recall, for all possible thresholds.

1. Which part of this curve corresponds to T = 0.9?
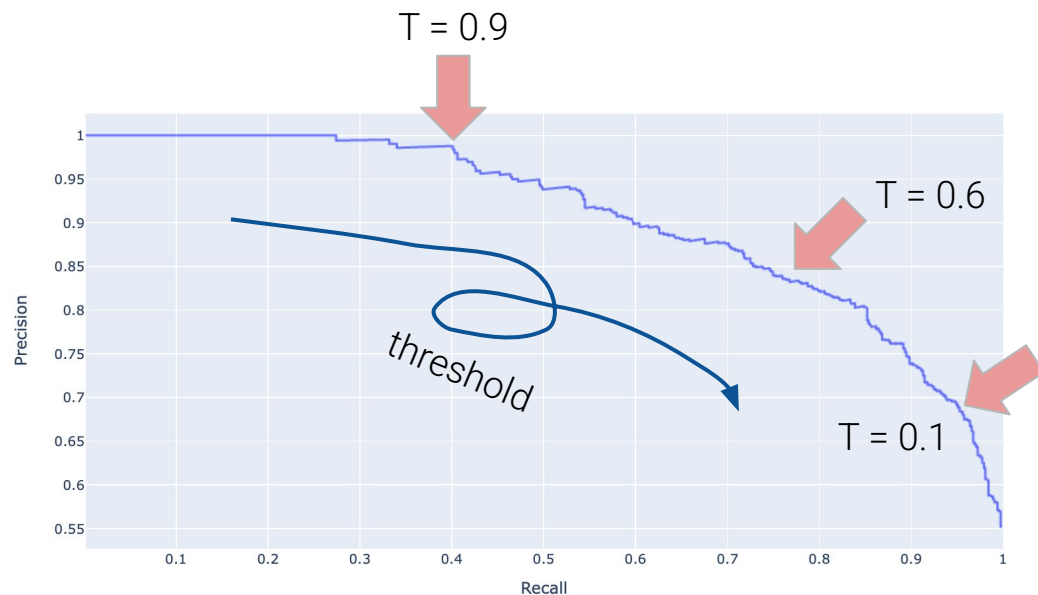2. Which part of this curve corresponds to T = 0.1?

# Precision-recall curves

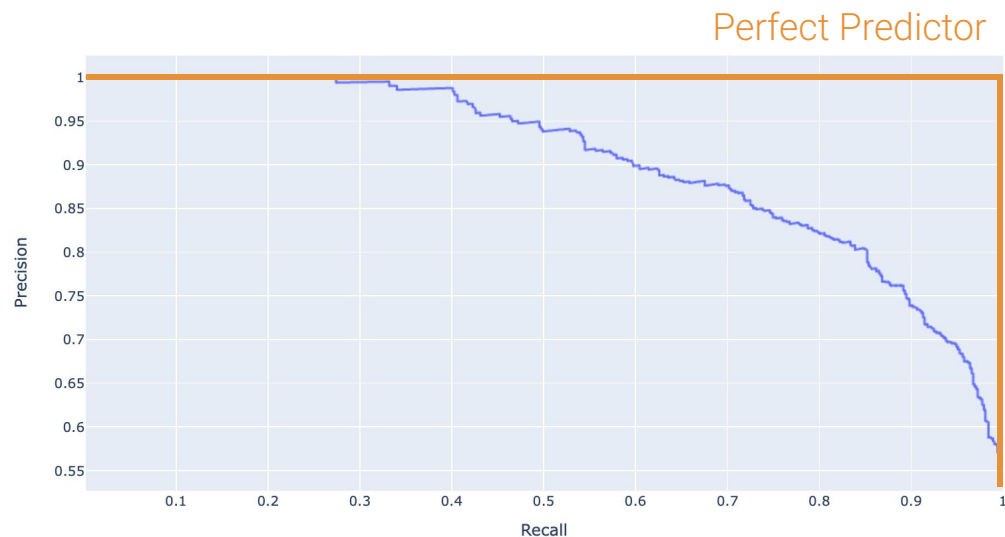We can also plot precision vs. recall, for all possible thresholds.

Answer:

- Threshold decreases from the top left to the bottom right.
- In the notebook, there's an interactive version of this plot.

# Precision-recall curves

The "perfect classifier" is one with precision of 1 and recall of 1.

- We want our PR curve to be as close to the "top right" of this graph as possible.
- One way to compare our model is to compute its **area under curve (AUC)**.
  - The area under the "optimal PR curve" is 1.
  - More commonly, we look at the area under ROC curve.



Perfect Predictor

# Logistic Regression II

Content credit: Lisa Yan, Suraj Rampure, Ani Adhikari, Josh Hug, Joseph Gonzalez