

Logistic Regression

Dr Mehrdad Ghaziasgar

Senior Lecturer in Computer Science
Department of Computer Science
University of the Western Cape
mghaziasgar@uwc.ac.za



UNIVERSITY of the
WESTERN CAPE

Content

- Part 1:
 - Classification With Logistic Regression
 - Model Evaluation For Classification
- Part 2:
 - Overfitting and Underfitting
- Part 3:
 - Practical Issues With Classification

Content - Part 1

- Classification With Logistic Regression
 - Model Representation
 - Obtaining the Decision Boundary
 - Cost Function
 - Cost Function Minimization For Classification
- Model Evaluation For Classification
 - Accuracy
 - Precision and Recall
 - F1 score

Classification With Logistic Regression

Model Representation

- Classification With Logistic Regression
 - Model Representation
 - Obtaining the Decision Boundary
 - Cost Function
 - Cost Function Minimization For Classification
 - Multi-Class Classification
- Model Evaluation For Classification
 - Accuracy
 - Precision and Recall
 - F1 score

Classification With Logistic Regression

- Example classification problems:
 - Credit-riskiness
 - Not Risky – Risky $y \in \{0,1\}$
 - Fraud detection
 - Not Fraudulent – Fraudulent $y \in \{0,1\}$
 - Sentiment detection / analysis
 - Negative – Positive $y \in \{0,1\}$
 - Neutral – Positive – Negative $y \in \{0,1,2\}$
 - Facial expression recognition
 - Neutral – Positive – Negative $y \in \{0,1,2\}$
 - Neutral – Happy – Sad – Surprise – Disgust – Angry $y \in \{0,1,2,3,4,5\}$
- y now takes on a set of specific (discrete) values

Classification With Logistic Regression

- With regression problems, we were trying to draw a line to **FIT** the data tightly
- With classification problems, we're now trying to draw a line to **SEPARATE** the data nicely:
 - E.g. For a two-class problem: Separate data points with class $y = 1$ from those with class $y = 0$
 - E.g. For a five-class problem: Separate data points of classes $y = 0$, $y = 1$, $y = 2$, $y = 3$ and $y = 4$ from each other
 - Etc.

Classification With Logistic Regression

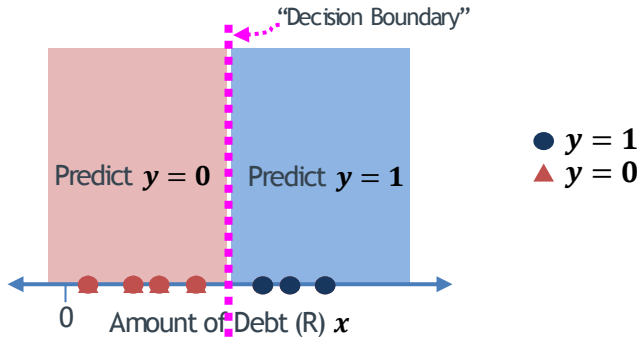
- With classification problems, we're now trying to draw a line to **SEPARATE** the data nicely
- E.g. The problem of predicting the credit riskiness (y) of a person based on one feature: amount of debt owed by a person (R) (x)

Amount of Debt (R) (x)	Credit Risk (y)
98000	0
549000	1
215000	0
69000	0
950000	1
...	...

Training set of credit riskiness

Classification With Logistic Regression

- With classification problems, we're now trying to draw a line to **SEPARATE** the data nicely
- E.g. 2 The problem of predicting the credit riskiness (y) of a person based on two features: amount of debt (R) (x_1) and the amount of payment defaults (R) that a person has had (x_2)



Classification With Logistic Regression

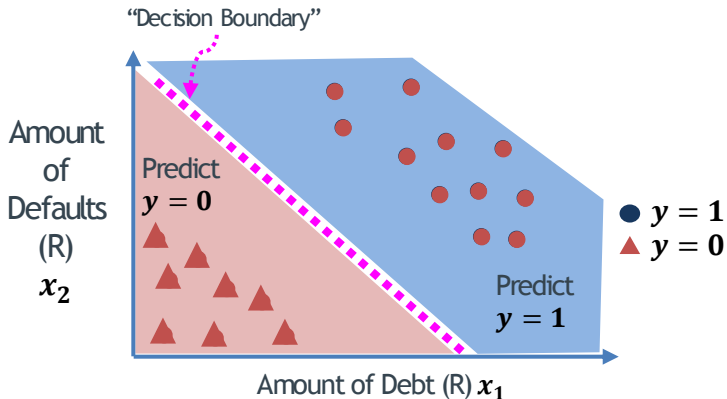
- With classification problems, we're now trying to draw a line to **SEPARATE** the data nicely
- E.g. 2 The problem of predicting the credit riskiness (y) of a person based on two features: amount of debt (R) (x_1) and the amount of payment defaults (R) that a person has had (x_2)

Amount of Debt (R) (x_1)	No. of Defaults (x_2)	Credit Risk (y)
98000	8150	0
549000	29000	1
215000	1750	0
69000	2000	0
950000	13000	1
...		...

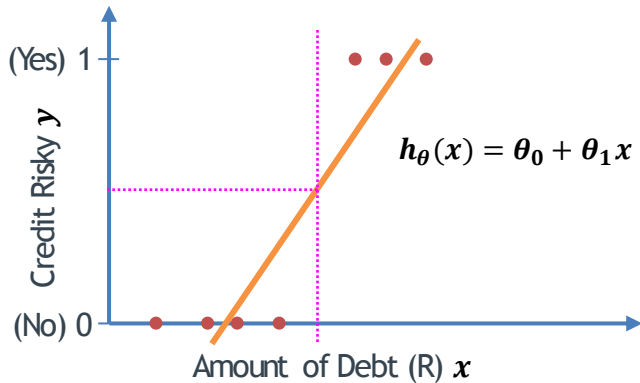
Training set of credit riskiness

Classification With Logistic Regression

- With classification problems, we're now trying to draw a line to **SEPARATE** the data nicely
- E.g. 2 The problem of predicting the credit riskiness (y) of a person based on two features: amount of debt (R) (x_1) and the amount of payment defaults (R) that a person has had (x_2)



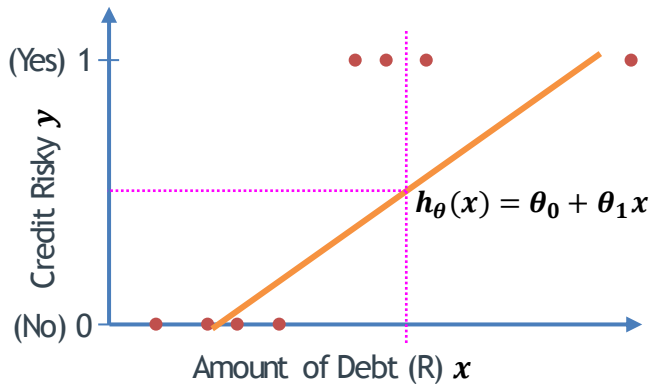
The Problem With Linear Regression



Possible Strategy(?):

- Use linear regression to fit a linear $h_{\theta}(x)$ to the data
- Then use a thresholding scheme to convert the linear scale into a discrete one e.g.:
 - If $h_{\theta}(x) < 0.5$ predict $y = 0$
 - If $h_{\theta}(x) \geq 0.5$ predict $y = 1$

The Problem With Linear Regression



Possible Strategy:

- Use linear regression to fit a linear $h_{\theta}(x)$ to the data
- Then use a thresholding scheme to convert the linear scale into a discrete one e.g.:
 - If $h_{\theta}(x) < 0.5$ predict $y = 0$
 - If $h_{\theta}(x) \geq 0.5$ predict $y = 1$

The Problem With Linear Regression

- Another problem with linear regression:
 - h_{θ} ranges from $-\infty$ to $+\infty$
 - Whereas our labels are now very specific and discrete i.e. $y \in \{0,1\}$
- What we want is a learning approach that:
 - Produces h_{θ} such that $0 \leq h_{\theta} \leq 1$
 - Focuses on **separating** (as opposed to **fitting**) data based on the labels
 - Separates the labels effectively regardless of noise in the data

Classification With Logistic Regression

- We will need a few elements to achieve classification:
 1. A (new) hypothesis $h_{\theta}(x)$ that can make predictions in the range $[0,1]$ given a feature vector x
 2. A (new) cost function $J(\theta)$ that gives us a increasingly large costs if the actual label of an example is 0 while the predicted label approaches 1 OR if the actual label of an example is 1 and the predicted label approaches 0
 3. A set of gradient descent update rules for θ using the new hypothesis and cost function so that we can learn the θ parameters

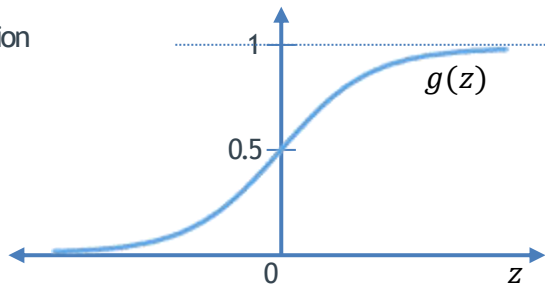
Model Representation

- We want a learning approach that:
 - Produces \mathbf{h}_θ such that $0 \leq \mathbf{h}_\theta \leq 1$
 - Focuses on **separating** (as opposed to **fitting**) data based on the labels
- We'll modify \mathbf{h}_θ from linear regression as follows:

$$\mathbf{h}_\theta(\mathbf{x}) = g(\boldsymbol{\theta}^T \mathbf{x}) \quad \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$$

- $g(z)$ is the “Sigmoid” or “Logistic” function

$$g(z) = \frac{1}{1 + e^{-z}}$$



Model Representation

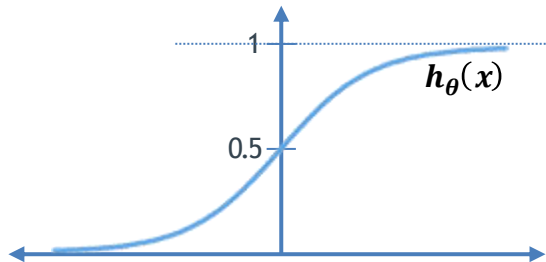
- Note that $g(z)$ is exactly in the range we want i.e. $0 \leq g(z) \leq 1$:
 - If z is large e.g. 100 then:
$$g(z) = \frac{1}{1 + e^{-z}}$$
 - If z is very small e.g. -100000 then:

- So now:

$$h_{\theta}(x) = g(\theta^T x)$$

- So:

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$



Model Representation - Deciphering (the new) h_θ

- The new h_θ can be thought of as a probability function
- It outputs the probability that a given set of features represent $y = 1$
- For a given x :
 - The closer h_θ is to 1, the more likely the class is $y = 1$
 - The closer h_θ is to 0, the more likely the class is $y = 0$
- E.g. In the credit risk example earlier:
 - Assume we have somehow determined the parameters $\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix}$
 - Given an arbitrary $x = \begin{bmatrix} x_0 = 1 \\ x_1 = \text{Amount of debt} \end{bmatrix}$
 - $h_\theta(x) = 0.2$ could mean 20% chance that client is a credit risk
 - $h_\theta(x) = 0.85$ could mean 85% chance that client is a credit risk

Model Representation - Deciphering h_{θ}

- In Mathematical notation:
 - $h_{\theta}(x) = P(y = 1|x; \theta)$
 - i.e. the probability that the class $y = 1$ given a specific example x and given the learned parameters θ
- As a by the way: Note that since we have only two classes $y = 0$ and $y = 1$, the sum of the probabilities of the two classes should be 1 i.e.

$$P(y = 1|x; \theta) + P(y = 0|x; \theta) = 1$$

$$P(y = 0|x; \theta) = 1 - P(y = 1|x; \theta)$$

$$P(y = 0|x; \theta) = 1 - h_{\theta}(x)$$

Model Representation - Deciphering h_θ

- As a by the way: Note that since we have only two classes $y = 0$ and $y = 1$, the sum of the probabilities of the two classes should be 1 i.e.

$$P(y = 0|x; \theta) = 1 - h_\theta(x)$$

- E.g. In the credit risk example earlier:

$$\text{Given an arbitrary } x = \begin{bmatrix} x_0 = 1 \\ x_1 = \text{Amount of debt} \end{bmatrix}$$

- If $h_\theta(x) = P(y = 1|x; \theta) = 0.8$ then the probability of the client being a credit risk is 80%
- So the probability of the same client NOT being a credit risk is:

$$P(y = 0|x; \theta) = 1 - 0.8 = 0.2 \text{ i.e. } 20\%$$

Classification With Logistic Regression

Obtaining the Decision Boundary

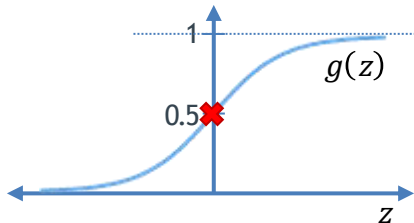
- Classification With Logistic Regression
 - Model Representation
 - Obtaining the Decision Boundary
 - Cost Function
 - Cost Function Minimization For Classification
 - Multi-Class Classification
- Model Evaluation For Classification
 - Accuracy
 - Precision and Recall
 - F1 score

Obtaining the Decision Boundary

- Here's what we have so far for logistic regression:

$$\mathbf{h}_{\theta}(\mathbf{x}) = g(\boldsymbol{\theta}^T \mathbf{x}) \quad g(z) = \frac{1}{1 + e^{-z}}$$

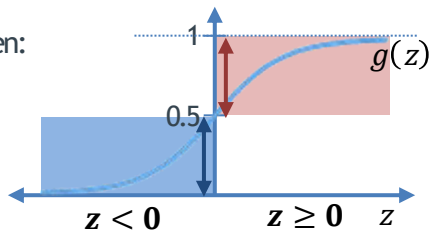
- We know that $\mathbf{h}_{\theta}(\mathbf{x})$ is the probability that $\mathbf{y} = 1$
- So we can make predictions by computing $\mathbf{h}_{\theta}(\mathbf{x})$ and then:
 - Predict $\mathbf{y} = 1$ whenever $\mathbf{h}_{\theta}(\mathbf{x}) \geq 0.5$
 - Predict $\mathbf{y} = 0$ whenever $\mathbf{h}_{\theta}(\mathbf{x}) < 0.5$



Obtaining the Decision Boundary

- So we can make predictions by computing $h_{\theta}(x)$ and then:

- Predict $y = 1$ whenever $h_{\theta}(x) \geq 0.5$ so:
- Predict $y = 1$ whenever $g(\theta^T x) \geq 0.5$
- But $g(z) \geq 0.5$ whenever $z \geq 0$
- So $g(\theta^T x) \geq 0.5$ whenever $\theta^T x \geq 0$ so finally:
- Predict $y = 1$ whenever $\theta^T x \geq 0$



- Similarly:

- Predict $y = 0$ whenever $h_{\theta}(x) < 0.5$ so:
- Predict $y = 0$ whenever $g(\theta^T x) < 0.5$
- But $g(z) < 0.5$ whenever $z < 0$
- So $g(\theta^T x) < 0.5$ whenever $\theta^T x < 0$ so finally:
- Predict $y = 0$ whenever $\theta^T x < 0$

Obtaining the Decision Boundary

- So:
 - Predict $y = 1$ whenever $\theta^T x \geq 0$
 - Predict $y = 0$ whenever $\theta^T x < 0$
- Now let's try to decipher:
 - what $\theta^T x$ represents
 - what $\theta^T x \geq 0$ and $\theta^T x < 0$ mean

Obtaining the Decision Boundary

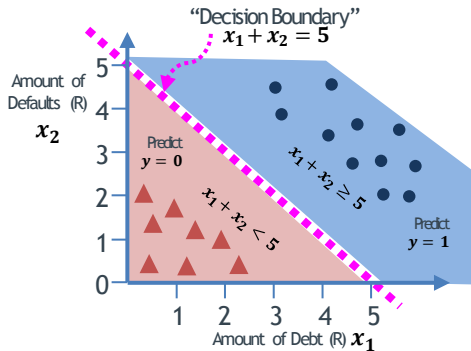
- Assume we've obtained $\theta = \begin{bmatrix} \theta_0 = -5 \\ \theta_1 = 1 \\ \theta_2 = 1 \end{bmatrix}$
- So $\theta^T x = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 = -5 + x_1 + x_2$
- And we now have:
- Predict $y = 1$ whenever $\theta^T x \geq 0$

$$-5 + x_1 + x_2 \geq 0$$

$$x_1 + x_2 \geq 5$$
- Predict $y = 0$ whenever $\theta^T x < 0$

$$-5 + x_1 + x_2 < 0$$

$$x_1 + x_2 \leq 5$$



- The “decision boundary” is given by:
- $\theta^T x = 0$ which is $-5 + x_1 + x_2 = 0$

$$x_1 + x_2 = 5$$
- Note that the decision boundary is defined by parameters θ only

Obtaining the Decision Boundary

- So to recap:
 - We obtain θ somehow (later - minimizing a cost function)
 - This θ defines (in logistic regression) a line that separates the two classes $y = 0$ and $y = 1$ called the “decision boundary”
 - Given any sample x , computing $h_{\theta}(x) = g(\theta^T x)$ tells us on which side of the decision boundary (or possibly on it) this sample falls
 - We use this ($h_{\theta}(x)$) to then determine whether x belongs to $y = 0$ or $y = 1$

Non-Linear Decision Boundaries

- In Linear Regression we were able to add non-linear features in order to **FIT** non-linear data more tightly:
 - Combinations of features
 - Higher-order features
- In Logistic Regression, we can add non-linear features in order to **SEPARATE** non-linear data (of different classes) more effectively:
 - Technique is exactly the same
 - The goal is different: to **separate** data classes
 - The hypothesis is different to that of Linear Regression

Non-Linear Decision Boundaries

- $h_{\theta}(x) = g(\theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$

- Assume we've obtained $\theta = \begin{bmatrix} \theta_0 = -4 \\ \theta_1 = 0 \\ \theta_2 = 0 \\ \theta_3 = 1 \\ \theta_4 = 1 \end{bmatrix}$

- Then:

$$\theta^T x = -4 \cdot 1 + 0 \cdot x_1 + 0 \cdot x_2 + 1 \cdot x_1^2 + 1 \cdot x_2^2$$

$$\theta^T x = -4 + x_1^2 + x_2^2$$

- So the decision boundary is $\theta^T x = 0$

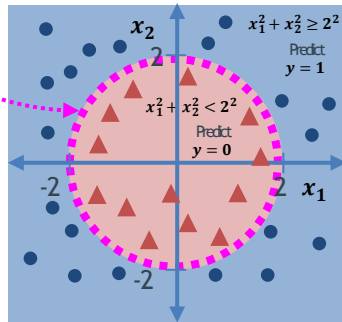
$$-4 + x_1^2 + x_2^2 = 0$$

$$x_1^2 + x_2^2 = 2^2$$

"Decision Boundary"

$$x_1^2 + x_2^2 = 2^2$$

● $y = 1$
▲ $y = 0$

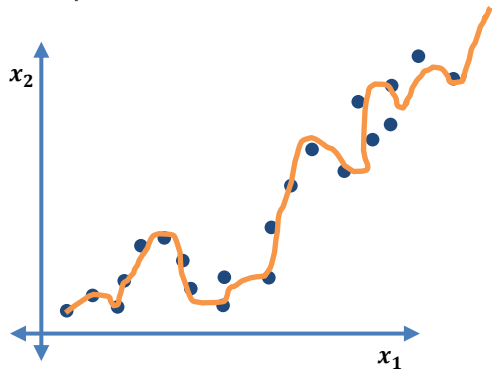


- Predict $y = 1$ whenever $\theta^T x \geq 0$
 $x_1^2 + x_2^2 \geq 2^2$

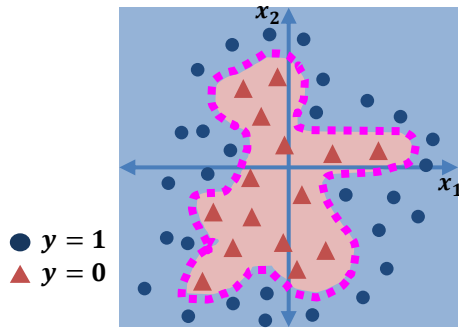
- Predict $y = 0$ whenever $\theta^T x < 0$
 $x_1^2 + x_2^2 < 2^2$

Non-Linear Decision Boundaries

- In Linear Regression, we were able to get virtually any complex non-linear fit-line to tightly **FIT** even the most complex data



- In Logistic Regression, using the same technique we can get virtually any kind of complex non-linear decision boundary by adding non-linear features to **SEPARATE** the data



Classification With Logistic Regression

Cost Function

- Classification With Logistic Regression
 - Model Representation
 - Obtaining the Decision Boundary
 - Cost Function
 - Cost Function Minimization For Classification
 - Multi-Class Classification
- Model Evaluation For Classification
 - Accuracy
 - Precision and Recall
 - F1 score

Cost Function - Can We Use $J(\theta)$ Of Linear Regression?

- We now have a (proven) valid hypothesis that can represent the separation of data based on two labels

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

- Now given a data set $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

Amount of Debt (R) (x_1)	No. of Defaults (x_2)	Credit Risk (y)
98000	8150	0
549000	29000	1
215000	1750	0
69000	2000	0
950000	13000	1
...		...

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \quad y^{(i)} \in \{0, 1\}$$

= 1

most effectively separates

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

- We need a cost function $J(\theta)$ that we can minimize

Training set of credit riskiness

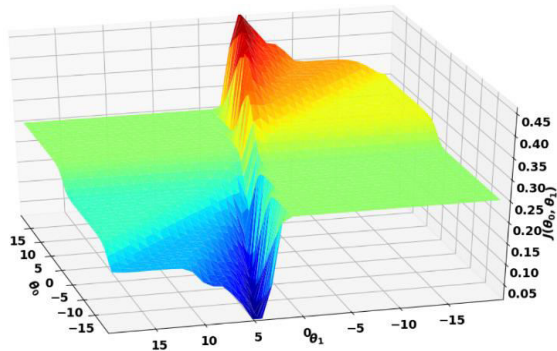
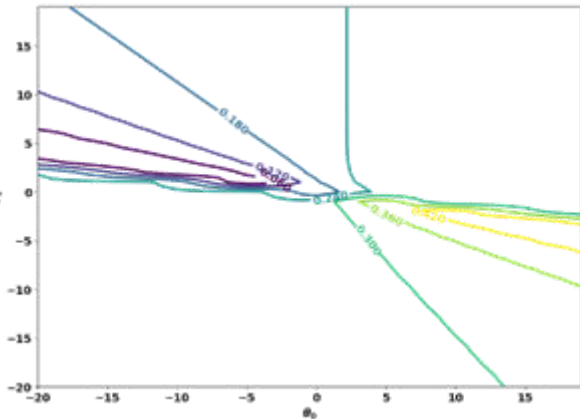
Cost Function - Can We Use $J(\theta)$ Of Linear Regression?

- The Linear Regression cost function was:

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

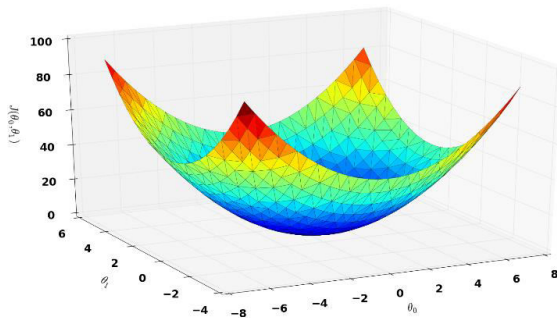
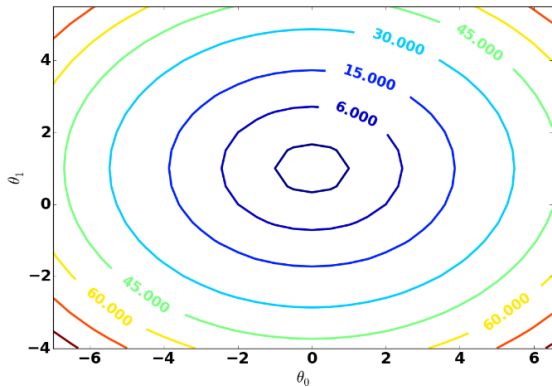
- If we use the Logistic Regression hypothesis $h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}}$ with this cost function, we get a complex non-convex cost function (feel free to confirm for yourself):
- Has many local minima:
 - Doesn't continuously decrease
 - Many regions where the function is completely flat
- We can't use gradient descent to minimize this
 - We'll end up in a local minimum

- E.g. For the Credit-Riskiness data set, using the least-squares cost function gives:



Cost Function

- We need a new/different cost function for Logistic Regression that will be convex (bowl-shaped) so we can minimize it



Cost Function

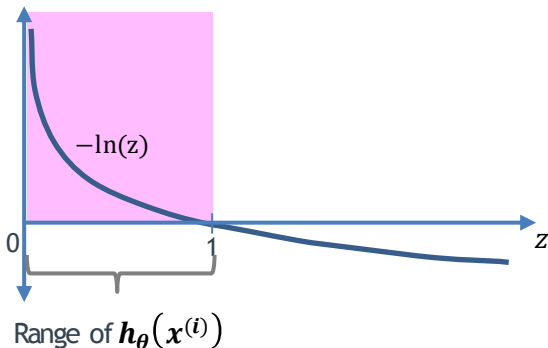
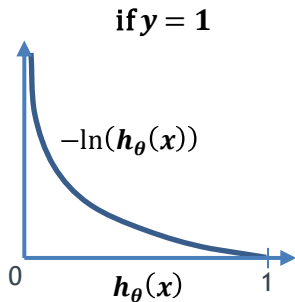
- We need a new/different cost function for Logistic Regression that will be convex (bowl-shaped) so we can minimize it
- Remember that $0 \leq h_{\theta}(x) \leq 1$
- The new cost function that we need should:
 - Have increasingly larger cost if the actual label is **if $y = 1$** but the prediction $h_{\theta}(x)$ approaches **0** OR the actual label is **if $y = 0$** but the prediction $h_{\theta}(x)$ approaches **1**
 - Have decreasing (minimum zero) cost if the prediction $h_{\theta}(x)$ approaches the actual label, whether it is $y = 0$ or $y = 1$.

Cost Function

- The following cost function has a convex shape with classification data:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(\mathbf{h}_{\theta}(\mathbf{x}^{(i)}), \mathbf{y}^{(i)})$$

$$\text{Cost}(\mathbf{h}_{\theta}(\mathbf{x}^{(i)}), \mathbf{y}^{(i)}) = \begin{cases} -\ln(\mathbf{h}_{\theta}(\mathbf{x}^{(i)})) & \text{if } \mathbf{y} = 1 \\ -\ln(1 - \mathbf{h}_{\theta}(\mathbf{x}^{(i)})) & \text{if } \mathbf{y} = 0 \end{cases}$$

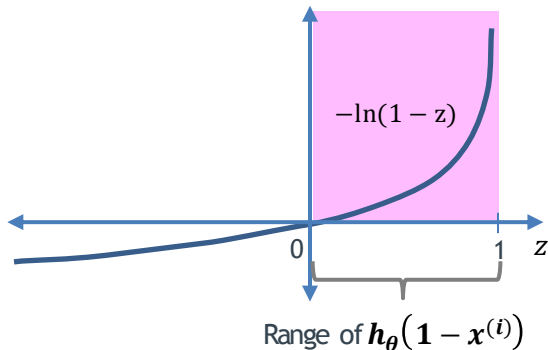
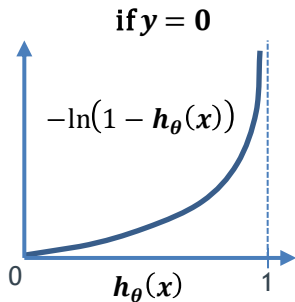


Cost Function

- The following cost function has a convex shape with classification data:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(\mathbf{h}_{\theta}(\mathbf{x}^{(i)}), \mathbf{y}^{(i)})$$

$$\text{Cost}(\mathbf{h}_{\theta}(\mathbf{x}^{(i)}), \mathbf{y}^{(i)}) = \begin{cases} -\ln(\mathbf{h}_{\theta}(\mathbf{x}^{(i)})) & \text{if } \mathbf{y} = 1 \\ -\ln(1 - \mathbf{h}_{\theta}(\mathbf{x}^{(i)})) & \text{if } \mathbf{y} = 0 \end{cases}$$

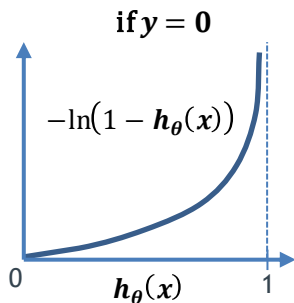
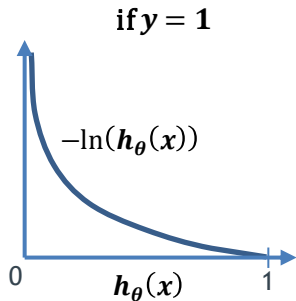


Cost Function

- The following cost function has a convex shape with classification data:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(\mathbf{h}_{\theta}(\mathbf{x}^{(i)}), \mathbf{y}^{(i)})$$

$$\text{Cost}(\mathbf{h}_{\theta}(\mathbf{x}^{(i)}), \mathbf{y}^{(i)}) = \begin{cases} -\ln(\mathbf{h}_{\theta}(\mathbf{x}^{(i)})) & \text{if } y = 1 \\ -\ln(1 - \mathbf{h}_{\theta}(\mathbf{x}^{(i)})) & \text{if } y = 0 \end{cases}$$



Cost Function

- We need a combined (not piece-wise defined) cost function to use with gradient descent
 - We need to combine the two parts of the function $\text{Cost}(\mathbf{h}_\theta(\mathbf{x}), \mathbf{y})$

$$\text{Cost}(\mathbf{h}_\theta(\mathbf{x}), \mathbf{y}) = \begin{cases} -\ln(\mathbf{h}_\theta(\mathbf{x})) & \text{if } \mathbf{y} = 1 \\ -\ln(1 - \mathbf{h}_\theta(\mathbf{x})) & \text{if } \mathbf{y} = 0 \end{cases}$$

- Remember that the labels have to be either $\mathbf{y} = 0$ or $\mathbf{y} = 1$
- If that's the case, one clever way of combining the two parts of $\text{Cost}(\mathbf{h}_\theta(\mathbf{x}), \mathbf{y})$ is:

$$\text{Cost}(\mathbf{h}_\theta(\mathbf{x}), \mathbf{y}) = -\mathbf{y} \ln(\mathbf{h}_\theta(\mathbf{x})) - (1 - \mathbf{y}) \ln(1 - \mathbf{h}_\theta(\mathbf{x}))$$

if $\mathbf{y} = 1$

$$\text{Cost}(\mathbf{h}_\theta(\mathbf{x}), \mathbf{y}) = -1 \cdot \ln(\mathbf{h}_\theta(\mathbf{x})) - (1 - 1) \ln(1 - \mathbf{h}_\theta(\mathbf{x}))$$

$$= -1 \cdot \ln(\mathbf{h}_\theta(\mathbf{x})) - (0) \cdot \ln(1 - \mathbf{h}_\theta(\mathbf{x}))$$

$$= -\ln(\mathbf{h}_\theta(\mathbf{x}))$$

Cost Function

- We need a combined (not piece-wise defined) cost function to use with gradient descent
 - We need to combine the two parts of the function $\text{Cost}(\mathbf{h}_\theta(\mathbf{x}), \mathbf{y})$

$$\text{Cost}(\mathbf{h}_\theta(\mathbf{x}), \mathbf{y}) = \begin{cases} -\ln(\mathbf{h}_\theta(\mathbf{x})) & \text{if } \mathbf{y} = 1 \\ -\ln(1 - \mathbf{h}_\theta(\mathbf{x})) & \text{if } \mathbf{y} = 0 \end{cases}$$

- Remember that the labels have to be either $\mathbf{y} = 0$ or $\mathbf{y} = 1$
- If that's the case, one clever way of combining the two parts of $\text{Cost}(\mathbf{h}_\theta(\mathbf{x}), \mathbf{y})$ is:

$$\text{Cost}(\mathbf{h}_\theta(\mathbf{x}), \mathbf{y}) = -\mathbf{y} \ln(\mathbf{h}_\theta(\mathbf{x})) - (1 - \mathbf{y}) \ln(1 - \mathbf{h}_\theta(\mathbf{x}))$$

if $\mathbf{y} = 0$

$$\begin{aligned} \text{Cost}(\mathbf{h}_\theta(\mathbf{x}), \mathbf{y}) &= -0 \cdot \ln(\mathbf{h}_\theta(\mathbf{x})) - (1 - 0) \ln(1 - \mathbf{h}_\theta(\mathbf{x})) \\ &= -0 \cdot \ln(\mathbf{h}_\theta(\mathbf{x})) - (1) \cdot \ln(1 - \mathbf{h}_\theta(\mathbf{x})) \\ &= -\ln(1 - \mathbf{h}_\theta(\mathbf{x})) \end{aligned}$$

Cost Function

- So with new combined version, we have: $J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(\mathbf{h}_{\theta}(\mathbf{x}^{(i)}), \mathbf{y}^{(i)})$

$$\begin{aligned}\text{Cost}(\mathbf{h}_{\theta}(\mathbf{x}^{(i)}), \mathbf{y}^{(i)}) &= -\mathbf{y}^{(i)} \ln(\mathbf{h}_{\theta}(\mathbf{x}^{(i)})) - (1 - \mathbf{y}^{(i)}) \ln(1 - \mathbf{h}_{\theta}(\mathbf{x}^{(i)})) \\ &= -\left[\mathbf{y}^{(i)} \ln(\mathbf{h}_{\theta}(\mathbf{x}^{(i)})) + (1 - \mathbf{y}^{(i)}) \ln(1 - \mathbf{h}_{\theta}(\mathbf{x}^{(i)}))\right]\end{aligned}$$

- So: $J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left[\mathbf{y}^{(i)} \ln(\mathbf{h}_{\theta}(\mathbf{x}^{(i)})) + (1 - \mathbf{y}^{(i)}) \ln(1 - \mathbf{h}_{\theta}(\mathbf{x}^{(i)}))\right]$
- We can now get the optimal parameters θ by minimizing this cost function $J(\theta)$ using gradient descent (or any other optimization technique)

$$\underset{\theta}{\text{minimize}} \ J(\theta)$$

Cost Function

- We can now get the optimal parameters θ by minimizing this cost function $J(\theta)$ using gradient descent (or any other optimization technique)

$$\underset{\theta}{\text{minimize}} \ J(\theta)$$

- We can then use the optimal θ to make predictions with the hypothesis $h_{\theta}(x)$:

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

Predict $y = 1$ whenever $h_{\theta}(x) \geq 0.5$

Predict $y = 0$ whenever $h_{\theta}(x) < 0.5$

Classification With Logistic Regression

Cost Function Minimization For Classification

- Classification With Logistic Regression
 - Model Representation
 - Obtaining the Decision Boundary
 - Cost Function
 - Cost Function Minimization For Classification
 - Multi-Class Classification
- Model Evaluation For Classification
 - Accuracy
 - Precision and Recall
 - F1 score

Cost Function Minimization For Classification

- With a cost function defined, we just need to minimize it to get the optimal θ

$$\underset{\theta}{\text{minimize}} \ J(\theta)$$

- There are a few ways to minimize the cost function (or any function):
 - Gradient Descent (we've already seen this)
 - Advanced optimization algorithms (we won't look at HOW these work - look it up):
 - Conjugate gradient descent
 - BFGS algorithm
 - L-BFGS algorithm
 - Nelder-Mead algorithm
 - Powell algorithm
 - Etc. (there are many)

Gradient Descent For Classification

- Generic gradient descent update equation (as seen before):

Repeat until convergence:

Update all θ_j simultaneously:

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

?

- The hypothesis and cost function are different to that of Linear Regression:

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

- Cost function:

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \ln(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \ln(1 - h_{\theta}(x^{(i)})) \right]$$

Gradient Descent for Classification

- For any parameter θ_j , the gradient is given by:

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Repeat until convergence:

Update all θ simultaneously:

$$\theta_0 \leftarrow \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_0^{(i)}$$

$$\theta_1 \leftarrow \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_1^{(i)}$$

...

$$\theta_n \leftarrow \theta_n - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_n^{(i)}$$

- Algorithm appears to be identical to gradient descent with Linear Regression
 - It is NOT! The hypothesis here is that of logistic regression
- The tuning of the learning rate follows a similar technique
- Feature scaling should be applied in the same way as explained before to achieve faster convergence
- Shouldn't use a for loop: Vectorize the approach

Advanced Optimization Algorithms For Classification

- Advanced optimization algorithms (we won't look at HOW these work - look it up):
 - Conjugate gradient descent
 - BFGS algorithm
 - L-BFGS algorithm
 - Nelder-Mead algorithm
 - Powell algorithm
 - Etc. (there are many)
- These algorithms take in the definition of the cost function and possibly the gradient values
 - They minimize the function using advanced techniques
 - No need to worry about learning rate
 - Work much faster than gradient descent
 - Implemented in SciPy

Advanced Optimization Algorithms For Classification

- Implemented in SciPy's optimize library:

```
from scipy.optimize import minimize
```

```
theta = np.zeros(X.shape[1]) #Initialize all thetas to zeros
```

```
result = minimize(costFunc, theta, args=(X,y), method='TNC', jac=gradientFunc,  
options={'maxiter' : 400, 'disp': True})
```

```
print(result.x)
```

- costFunc is a function that returns the cost as per:

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \ln(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \ln(1 - h_{\theta}(x^{(i)})) \right]$$

- gradientFunc is a function that returns a list/vector of the gradients of all the θ_j as per:

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Advanced Optimization Algorithms For Classification

- Implemented in SciPy's optimize library:

```
from scipy.optimize import minimize
```

```
theta = np.zeros(X.shape[1]) #Initialize all thetas to zeros
```

```
result = minimize(costFunc, theta, args=(X,y), method='TNC', jac=gradientFunc,  
options={'maxiter' : 400, 'disp': True})
```

```
print(result.x)
```

- result.x will contain the optimal θ values
- Note that X must contain an extra column of zeros representing feature x_0
- Predictions can then be made by applying $\mathbf{h}_{\theta}(\mathbf{x}) = g(\theta^T \mathbf{x})$

Advanced Optimization Algorithms For Classification

- Implemented in sklearn specifically for Logistic Regression:

```
from sklearn.linear_model import LogisticRegression
```

```
clf = LogisticRegression(random_state=0, solver='lbfgs', multi_class='auto')
```

```
clf.fit(X, y)
```

- Note that X should NOT contain the extra column of 1s for feature x_0 .
- Predictions can then be made seamlessly:

```
predictionvalue = clf.predict(x_new)
```

```
print(predictionvalue)
```

Classification With Logistic Regression

Multi-Class Classification

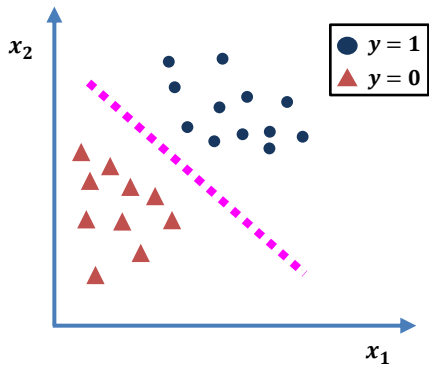
- Classification With Logistic Regression
 - Model Representation
 - Obtaining the Decision Boundary
 - Cost Function
 - Cost Function Minimization For Classification
 - Multi-Class Classification
- Model Evaluation For Classification
 - Accuracy
 - Precision and Recall
 - F1 score

Multi-Class Classification

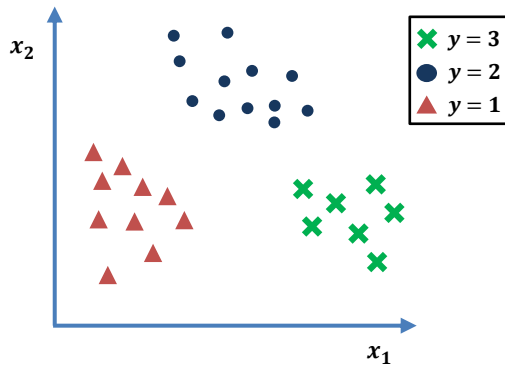
- Two broad types of classification
- **Binary classification;** $y \in \{0,1\}$
 - Two classes E.g.
 - Credit-riskiness
 - Not Risky – Risky
 - Fraud detection
 - Non-fraud – Fraud
 - Sentiment detection / analysis
 - Negative – Positive
 - Facial expression recognition
 - Emotive – Non-Emotive
- **Multi-class classification;** $y \in \{0,1,2 \dots, k\}$
 - Three or more classes E.g.
 - Credit-riskiness level
 - 0 Risk – 1 Risk – 2 Risk – 3 Risk $y \in \{0,1,2,3\}$
 - Sentiment detection / analysis
 - Neutral – Positive – Negative $y \in \{0,1,2\}$
 - Facial expression recognition
 - Neutral – Positive – Negative $y \in \{0,1,2\}$
 - Neutral – Happy – Sad – Surprise – Disgust – Angry $y \in \{0,1,2,3,4,5\}$

Multi-Class Classification

- Two broad types of classification
- Binary classification; $y \in \{0,1\}$



- Multi-class classification; $y \in \{1,2,3\}$

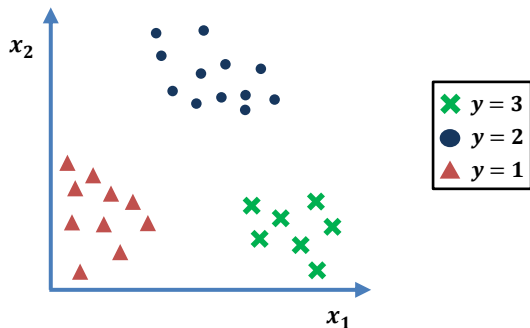


Multi-Class Classification

- Two main methods of achieving multi-class classification:
 - One-vs-Rest
 - One-vs-One

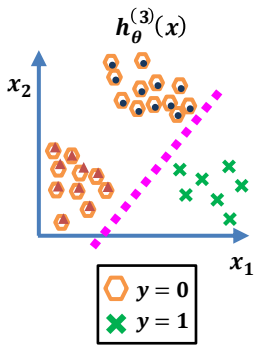
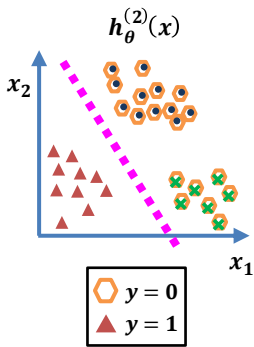
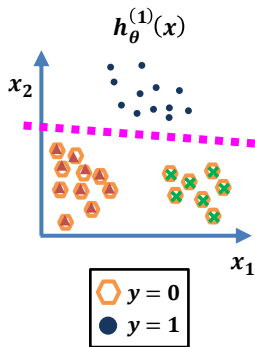
Multi-Class Classification - One-Vs-Rest

- One-vs-Rest:
- Decompose the K classes into K binary classification problems
- Create K separate classifiers $h_{\theta}^{(i)}(x)$ where $i = 1, 2, \dots, K$.
- In each classifier, set one of the K classes as $y = 1$ and combine all other classes into a single class $y = 0$
- E.g. for a 3-class problem:



Multi-Class Classification - One-Vs-Rest

- One-vs-Rest:
- Decompose the K classes into K binary classification problems
- Create K separate classifiers $h_{\theta}^{(i)}(x)$ where $i = 1, 2, \dots, K$.
- In each classifier, set one of the K classes as $y = 1$ and combine all other classes into a single class $y = 0$
- E.g. for a 3-class problem:



Multi-Class Classification - One-Vs-Rest

- One-vs-Rest:
- Decompose the K classes into K binary classification problems
- Create K separate classifiers $h_{\theta}^{(i)}(x)$ where $i = 1, 2, \dots, K$.
- In each classifier, set one of the K classes as $y = 1$ and combine all other classes into a single class $y = 0$
- Remember that the prediction of $h_{\theta}^{(i)}(x)$ is actually a probability; in this case, the probability that $y = i$:

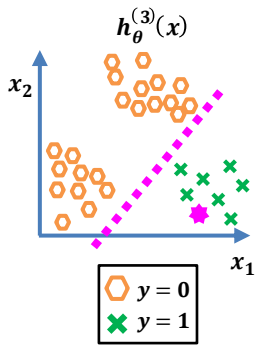
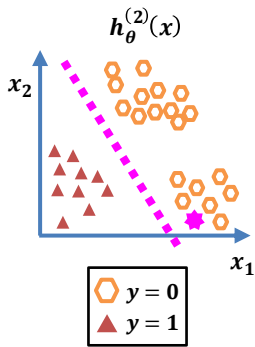
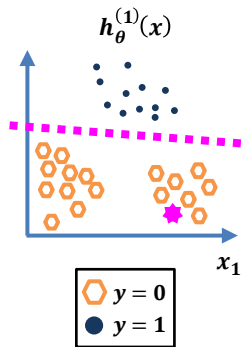
$$h_{\theta}^{(i)}(x) = P(y = 1|x; \theta)$$

- Then, given an unknown sample x , pass it to all $h_{\theta}^{(i)}(x)$ classifiers: classifier i with the highest probability means that x most likely belongs to class i
- Mathematically:

$$\max_i h_{\theta}^{(i)}(x)$$

Multi-Class Classification - One-Vs-Rest

- E.g. With the 3-class problem: assume we have a new x that actually belongs to class \times i.e. $y = 3$
- Pass x to $h_{\theta}^{(1)}(x)$, $h_{\theta}^{(2)}(x)$ and $h_{\theta}^{(3)}(x)$
- E.g. $h_{\theta}^{(1)}(x)$ will predict 0.01; $h_{\theta}^{(2)}(x)$ will predict 0.15; and $h_{\theta}^{(3)}(x)$ will predict 0.97. Therefore the predicted class is taken to be $y = 3$



Multi-Class Classification - One-Vs-One

- One-vs-One:
- Create a series of binary classifiers $h_{\theta}^{(i,j)}(x)$, which deal with only a pair of classes i and j .
- Remember that the prediction of $h_{\theta}^{(i,j)}(x)$ is actually a probability; in this case, the probability that $y = i$:

$$h_{\theta}^{(i)}(x) = P(y = 1|x; \theta)$$

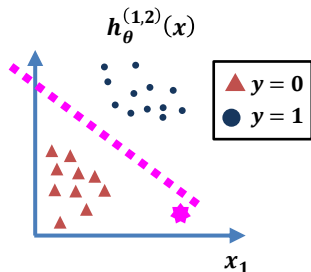
- But we can also infer the probability of class $y = 0$:

$$P(y = 0|x; \theta) = 1 - h_{\theta}^{(i)}(x)$$

- In this technique, we keep track of the total probability of each class:
 - The class with the max probability wins

Multi-Class Classification - One-Vs-One

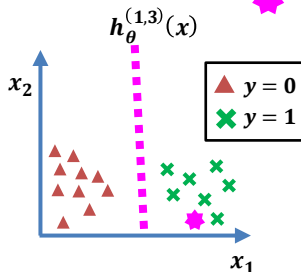
- E.g. With the 3-class problem: assume we have a new x that actually belongs to class \times i.e. $y = 3$



$h_{\theta}^{(i)}(x) = 0.48$ SO:

● $P(y = 1|x; \theta) = 0.48$

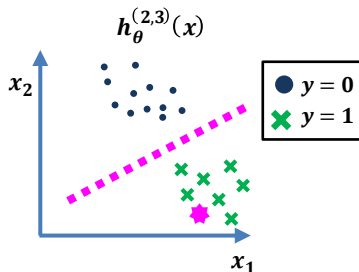
▲ $P(y = 0|x; \theta) = 0.52$



$h_{\theta}^{(i)}(x) = 0.97$ SO:

✕ $P(y = 1|x; \theta) = 0.97$

▲ $P(y = 0|x; \theta) = 0.03$



$h_{\theta}^{(i)}(x) = 0.99$ SO:

✕ $P(y = 1|x; \theta) = 0.99$

● $P(y = 0|x; \theta) = 0.01$

Totals:

● $= 0.48 + 0.01$
 $= 0.49$

▲ $= 0.52 + 0.03$
 $= 0.55$

✕ $= 0.97 + 0.99$
 $= 1.96$

Multi-Class Classification - One-Vs-One

- One-vs-One:
- Number of classifiers trained:
 - 3 classes: 3 classifiers i.e. classes 1-2, classes 1-3, classes 2-3
 - 4 classes: 6 classifiers i.e. classes 1-2, classes 1-3, classes 1-4, classes 2-3, classes 2-4, classes 3-4
 - K classes: $\frac{1}{2} K(K - 1)$

Multi-Class Classification

- One-vs-Rest:
 - Computationally cheaper and simpler
 - Fewer classifiers are trained
 - Fewer classifiers are used to make predictions
 - More sensitive to unbalanced classes (find out why)
 - Can be less accurate
- One-vs-One:
 - Computationally more expensive
 - Many more classifiers are trained
 - Many more classifiers are used to make predictions
 - Less sensitive to unbalanced classes (find out why)
 - Can be more accurate

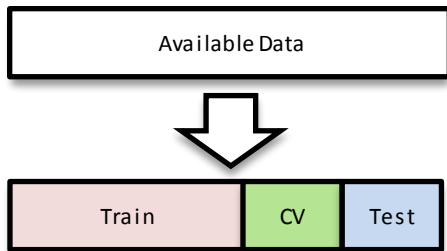
Model Evaluation For Classification

Accuracy

- Classification With Logistic Regression
 - Model Representation
 - Obtaining the Decision Boundary
 - Cost Function
 - Cost Function Minimization For Classification
 - Multi-Class Classification
- Model Evaluation For Classification
 - Accuracy
 - Precision and Recall
 - F1 score

Introduction

- When evaluating/comparing models, all the techniques described under Linear Regression apply i.e.
 - Divide up data into Train-CV-Test OR:
 - k -fold Cross-validation
- Question: what metrics do we use for classification?
- The metrics change for classification:
 - Accuracy
 - Precision
 - Recall
 - F1 score
- Running example: credit-riskiness
 - Assume we've trained a model $\mathbf{h}_{\theta}(\mathbf{x}^{(i)})$ for credit-riskiness on the Train set $\mathbf{x}_{\text{Train}}^{(i)}$
 - We need to evaluate the model on the Test set $\mathbf{x}_{\text{Test}}^{(i)}$



Accuracy

- Accuracy: the most basic metric
- Apply the test set samples $\mathbf{x}_{\text{Test}}^{(i)}$ to the trained model $\mathbf{h}_{\theta}(\mathbf{x}^{(i)})$ to get prediction labels for each sample
- The accuracy is the proportion/percentage of test samples for which the actual label $\mathbf{y}_{\text{Test}}^{(i)}$ matched the predicted label $\mathbf{h}_{\theta}(\mathbf{x}_{\text{Test}}^{(i)})$:

$$\text{Accuracy} = \frac{\text{Total Correct}}{\text{Total Test Samples}} \times 100$$

Model Evaluation For Classification

Precision and Recall

- Classification With Logistic Regression
 - Model Representation
 - Obtaining the Decision Boundary
 - Cost Function
 - Cost Function Minimization For Classification
 - Multi-Class Classification
- Model Evaluation For Classification
 - Accuracy
 - Precision and Recall
 - F1 score

Precision and Recall - The Problem

- Skew classes:
 - A case in which one of the classes to be predicted is scarce in the data set
 - In such data sets, we may not have many examples of the one of the classes
- E.g.
 - Fraud: there may be far fewer examples of fraud cases than legitimate cases
 - Cancer: there may be far fewer examples of cancerous cases than non-cancerous cases
- In such cases, one class (usually the class we're interested in) may only make up a small percentage of the examples e.g. only 0.5% of the examples
- The classifier may become very good at detecting negative cases, but not positive cases
- Using accuracy will not be a good reflection of the classifier performance

Precision and Recall - The Problem

- Skew classes:
 - A case in which one of the classes to be predicted is scarce in the data set
 - In such data sets, we may not have many examples of the one of the classes
- E.g. For the credit-riskiness data set, assume we have a test set of 1500 negative (non-risky $y = 0$) examples, and only 50 positive (risky $y = 1$) examples
- If our classifier correctly predicts 1300 negative samples, but only 10 of the positive examples, accuracy is?
 - $\frac{1300+10}{1500+50} \times 100 = 84.5\%$
 - But is this classifier really good at predicting credit-riskiness?? (NO!)

Precision and Recall - The Problem

- Skew classes:
 - A case in which one of the classes to be predicted is scarce in the data set
 - In such data sets, we may not have many examples of the one of the classes
- E.g. For the credit-riskiness data set, assume we have a test set of 1500 negative (non-risky **$y=0$**) examples, and only 50 positive (risky **$y=1$**) examples
- In fact, considering the skewness of the classes, we can use a dummy function to get an even better accuracy:

```
def predict(X):  
    return 0 #Just ignore X. Assume the class is 0
```

- Using this function to “predict” will give us an accuracy of?
 - $\frac{1500}{1500+50} \times 100 = 96.8\%$
- Accuracy is not a true reflection of underlying performance at all (in such cases)

Precision and Recall - The Problem

- We need some other way of evaluating our classifier that is not sensitive to the class distribution of examples - we're going to work towards it
- Four things that are important:
 - Of the cases that are actually **positive** (e.g. risky, fraudulent, cancerous etc.):
 1. how many are predicted **correctly** by the classifier?
 2. how many are predicted **wrongly** by the classifier?
 - Of the cases that are actually **negative** (e.g. reliable, not fraudulent, benign etc.):
 3. how many are predicted **correctly** by the classifier?
 4. how many are predicted **wrongly** by the classifier?
- We don't want someone who is risky, fraudulent, cancerous etc. to slip through the cracks
- We don't want to tell someone they are risky, fraudulent, cancerous etc. if they aren't

Precision and Recall - Confusion Matrix

- The four questions on the previous slide have specific “names” in ML literature:
 - Of the cases that are **actually positive**:
 - how many are **predicted correctly** by the classifier?
 - how many are **predicted wrongly** by the classifier?
 - Of the cases that are **actually negative**:
 - how many are **predicted correctly** by the classifier?
 - how many are **predicted wrongly** by the classifier?
- We can represent this information in a table called a “Confusion Matrix”

predicted		
actually	1.	2.
	4.	3.

Precision and Recall - Confusion Matrix

- The four questions on the previous slide have specific “names” in ML literature:
 - Of the cases that are **actually positive**:
 - how many are **predicted correctly** by the classifier?
 - how many are **predicted wrongly** by the classifier?
 - Of the cases that are **actually negative**:
 - how many are **predicted correctly** by the classifier?
 - how many are **predicted wrongly** by the classifier?
- In the label True/False **Positive/Negative**:
 - Positive/Negative** refers to the prediction
 - True/False refers to whether the prediction was right or wrong

		predicted	
		1	0
actually	1		
	0		

Precision and Recall - Confusion Matrix

- E.g. For credit riskiness, fill in the confusion matrix for:
 - 50 positive examples of which 10 were correctly predicted
 - 1500 negative examples of which 1300 were correctly predicted
- TP: FN:
- FP: TN:

		predicted	
		1	0
actually	1		
	0		

Precision and Recall - Confusion Matrix

- E.g. For credit riskiness, fill in the confusion matrix for:
 - 15 positive examples of which 7 were correctly predicted
 - 270 negative examples of which 200 were correctly predicted
- TP FN:
- FP TN:

		predicted	
		1	0
actually	1		
	0		

Precision and Recall - Confusion Matrix

- E.g. For cat detection, fill in the confusion matrix for:
 - 850 positive examples of which 420 were INcorrectly predicted
 - 650 negative examples of which 360 were INcorrectly predicted
- TP: FN:
- FP: TN:

		predicted	
		1	0
actually	1		
	0		

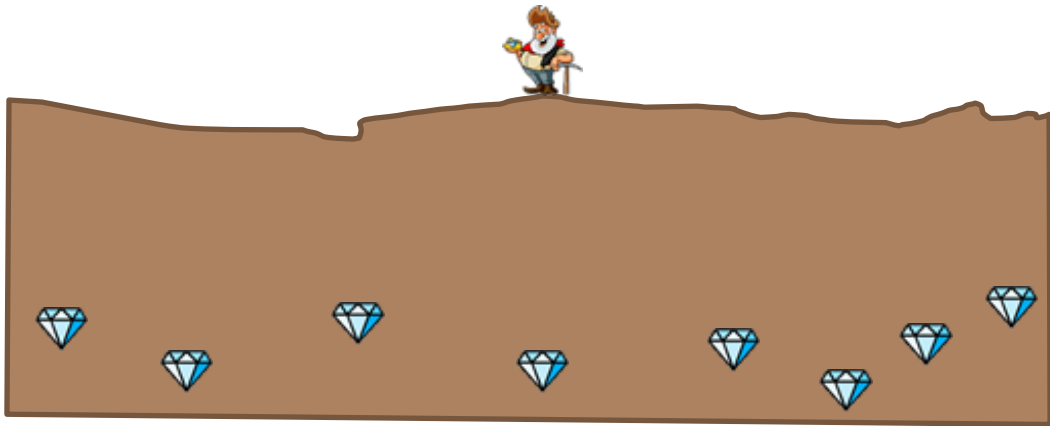
Precision and Recall - Confusion Matrix

- E.g. For cancer detection, fill in the confusion matrix for:
 - 160 positive examples of which 10 were Incorrectly predicted
 - 1750 negative examples of which 600 were correctly predicted
- TP: FN:
- FP: TN:

		predicted	
		1	0
actually	1		
	0		

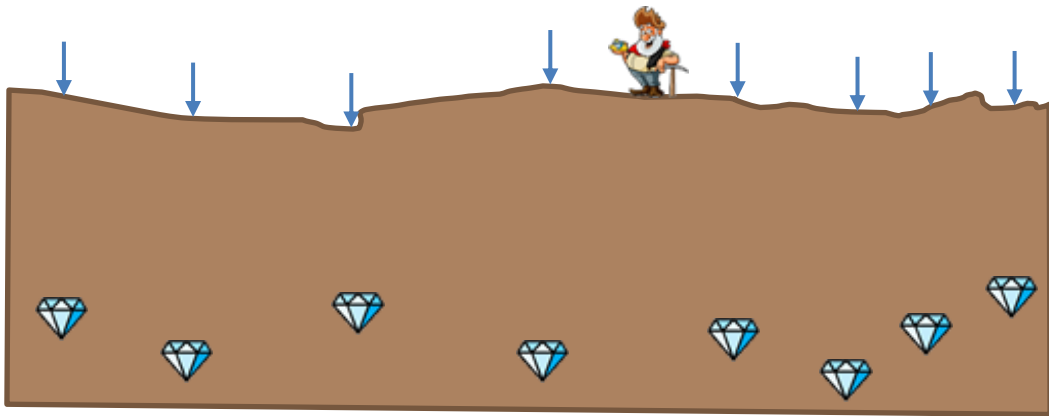
Precision and Recall - The Story

- You've employed a diamond expert to help you dig for diamonds
- The main question you want him to answer: where should you dig?
- When he says "dig here", do you trust him (precision)? How many of the diamonds will he find (recall)?



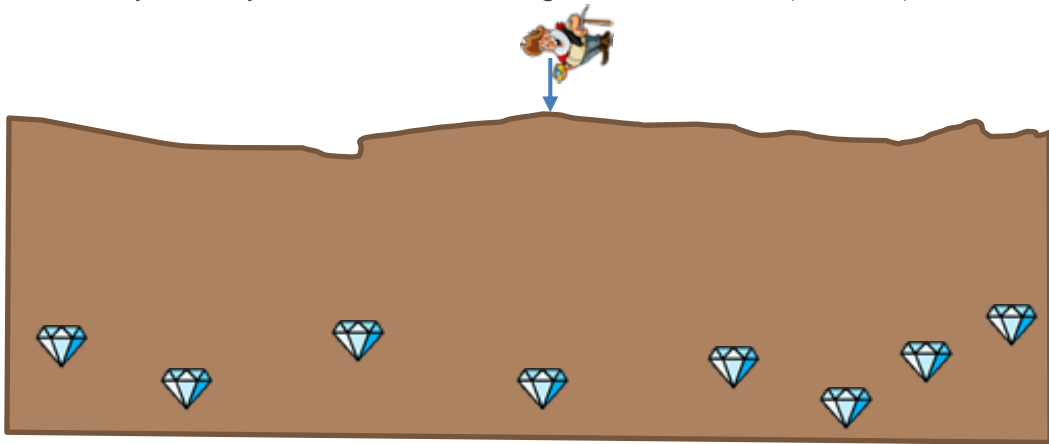
Precision and Recall - Interpretation

- Perfect Case:
 - Whenever he says “Diamond”, it definitely is... DIG! (High Precision)
 - He finds ALL the diamonds! 💎 (High Recall)




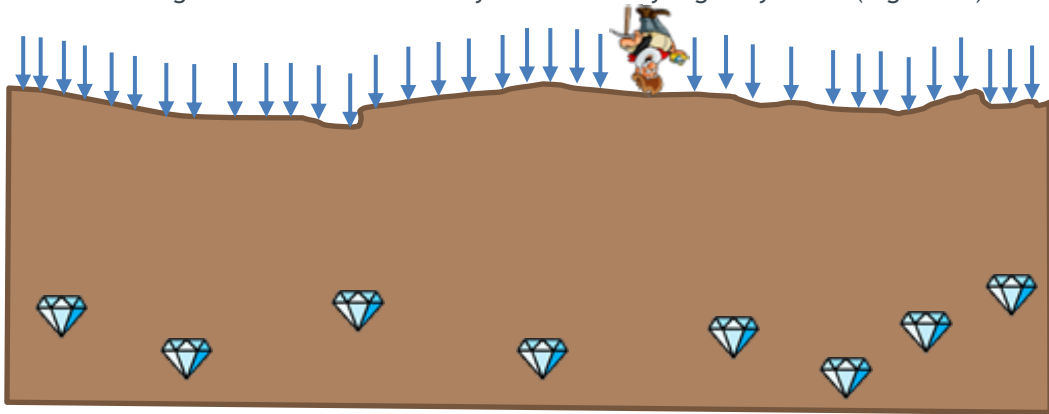
Precision and Recall - Interpretation

- Non-Ideal Case:
 - Whenever he says “Diamond”, it definitely is... DIG! (High Precision)
 - He only finds very few diamonds! 💎 Is this gonna be sustainable?? (Low Recall)



Precision and Recall - Interpretation

- Non-Ideal Case:
 - When he says “Diamond”, it is unlikely that it is... dig?? Rather ask your cat... (Low Precision)
 - He finds many/all the diamonds!  In essence, he’s just saying “DIG EVERYWHERE. You’re bound to get all the diamonds that way...” Can we really dig everywhere? (High Recall)



Precision

- For any classifier, there are two critical questions:
 1. **Precision:** How **reliable** / **trustworthy** / **precise** are the classifier's **predictions**?
 - Can we trust the classifier's predictions?
 - When the classifier throws an alarm (i.e. $y = 1$, risky, fraudulent, cancerous) how sure can we be that it actually is? We wouldn't want to react unnecessarily
 - When the classifier says all is fine (i.e. $y = 0$, not risky, not fraudulent, not cancerous) how sure can we be that it actually is?
 - Another way of saying this: of the cases that the classifier **predicts** as **positive**, what proportion are **predicted correctly**?

Precision

- For any classifier, there are two critical questions:
 - Precision:** How **reliable** / **trustworthy** / **precise** are the classifier's **predictions**?
 - Another way of saying this: of the cases that the classifier **predicts as positive**, what proportion are **predicted correctly**?
 - Cases **predicted correctly**:
 - Cases **predicted as positive**:
(all of the cases predicted as positive)

		predicted	
		1	0
actually	1	True Positives	False Negatives
	0	False Positives	True Negatives

$$\begin{aligned}
 \text{Precision} &= \frac{\text{True Positives}}{\text{Total Positive Predictions}} (\times 100) \\
 &= \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}} (\times 100)
 \end{aligned}$$

Recall

- For any classifier, there are two critical questions:
 2. **Recall:** How **comprehensive/thorough** is the classifier in finding cases that are **actually positive**?
 - How many of the actually positive cases can the classifier actually find?
 - How sure can we be that the classifier will find all of the positive cases?
 - How sure can we be that all of the positive cases (i.e. $y = 1$, risky, fraudulent, cancerous) won't just slip through undetected?
 - When the classifier identifies a bunch of cases as being positive, how sure can we be that these are all of them?
 - We wouldn't want to miss positive cases which need to be reacted on
 - Another way of saying this: of the cases that are **actually positive**, what proportion are **predicted correctly**?

Recall

- For any classifier, there are two critical questions:
 - Recall:** How **comprehensive/thorough** is the classifier in finding cases that are **actually positive**?
 - Another way of saying this: of the cases that are **actually positive**, what proportion are **predicted correctly**?
 - Cases **predicted correctly**:
 - Cases **actually positive**:
(all of the actually positive examples)

		predicted	
		1	0
actually	1	True Positives	False Negatives
	0	False Positives	True Negatives

$$\begin{aligned}
 \text{Recall} &= \frac{\text{True Positives}}{\text{Total Actual Positive Examples}} (\times 100) \\
 &= \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}} (\times 100)
 \end{aligned}$$

Precision and Recall - The Problem (Revisited)

- E.g. For the credit-riskiness data set, assume we have a test set of 1500 negative (non-risky $y = 0$) examples, and only 50 positive (risky $y = 1$) examples
- If our classifier correctly predicts 1300 negative samples, but only 10 of the positive examples, accuracy was:
 - $\frac{1300+10}{1500+50} \times 100 = 84.5\%$
- What about precision?
 - $\frac{TP}{TP+FP} = \frac{10}{10+200} \times 100 \approx 5\%$
- And recall?
 - $\frac{TP}{TP+FN} = \frac{10}{10+40} \times 100 = 20\%$
- Accuracy was deceiving, but not precision and recall. This classifier is awful

		predicted	
		1	0
actually	1		
	0		

Precision and Recall - The Problem (Revisited)

- E.g. For the credit-riskiness data set, assume we have a test set of 1500 negative (non-risky $y=0$) examples, and only 50 positive (risky $y=1$) examples

- What about the dummy function?

```
def predict(X):  
    return 0    #Just ignore X. Assume the class is 0
```

- Using this function to “predict” gave us an accuracy of:

- $$\frac{1500}{1500+50} \times 100 = 96.8\%$$

- What's the precision of the dummy function?

- $$\frac{TP}{TP+FP} = \frac{0}{0+0} \times 100 = [\text{NaN}] \sim 0\%$$

- What's the recall of the dummy function?

- $$\frac{TP}{TP+FN} = \frac{0}{0+50} \times 100 = 0\%$$

- Accuracy was deceiving, but not precision and recall

		predicted	
		1	0
actually	1		
	0		

Precision and Recall

- E.g. For cat detection:
 - 275 positive examples of which 200 were INcorrectly predicted
 - 260 negative examples of which 20 were INcorrectly predicted

- Accuracy:

-

- Precision:

-

- Recall:

-

		predicted	
		1	0
actually	1	75	200
	0	20	240

- Interpretation: (High Prec): When it says it's a cat, very good (79%) chance it is; (Low Rec) BUT it won't find a lot of these cats... a lot of the cats will simply go undetected (as non-cats)...

Precision and Recall

- E.g. For cancer detection:
 - 160 positive examples of which 10 were INcorrectly predicted
 - 1750 negative examples of which 600 were correctly predicted

- Accuracy:
 -

- Precision:
 -

- Recall:
 -

		predicted	
		1	0
actually	1	150	10
	0	1150	600

- Interpretation: (High Recall): It'll find almost all cancer cases and "treat" them; (Low Prec) BUT it'll also end up telling a lot of people who aren't cancerous that they are... and possible "treat" them too...

Precision and Recall - Interpretation

- Ideal Case:
 - High Precision - High Recall: Whenever it says it's a positive it most likely is; and it will successfully locate (almost) all positive cases (few/none will slip through)
- Non-Ideal Cases:
 - High Precision - Low Recall: Whenever it says it's a positive it most likely is; BUT it will locate very few (if any) positive cases (many will slip through)
 - High Recall - Low Precision: It will locate (almost) all positive cases (few/none will slip through); BUT when it says it's a positive, it will most likely not be one - we can't trust it

Model Evaluation For Classification

F1 Score

- Classification With Logistic Regression
 - Model Representation
 - Obtaining the Decision Boundary
 - Cost Function
 - Cost Function Minimization For Classification
 - Multi-Class Classification
- Model Evaluation For Classification
 - Accuracy
 - Precision and Recall
 - F1 score

F1 Score - The Problem

- We've now got two numbers (precision and recall)
 - They give us a very good understanding of our classifier's performance
 - BUT they are still TWO numbers
 - We need some way of combining the two numbers into one number
- E.g. For the credit-riskiness example: imagine we have three different hypotheses:
 1. $h_{\theta}(x) = g(\theta_0 x_0 + \theta_1 x_1)$
 2. $h_{\theta}(x) = g(\theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2)$
 3. $h_{\theta}(x) = g(\theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$
- Where x_1 is the No of Defaults and x_2 is the Total Debt Owed (R)
- How do these hypotheses compare??
 - One way: compute precision and recall

F1 Score - The Problem

- We get these results:

Approach	Precision (%)	Recall (%)
1	60	40
2	12	98
3	95	5

- Which model is best? We need some way of combining the two numbers into one number
- One possible way: taking the average??

F1 Score - The Problem

- We get these results:

Approach	Precision (%)	Recall (%)	Average (%)
1	60	40	50
2	12	98	55
3	95	5	50

- The average doesn't seem to be a good indicator for performance:
 - All these approaches appear to be almost exactly the same considering the average whereas:
 - Algorithm 1 is more balanced between precision and recall
 - Algorithm 2 simply says "Dig everywhere"
 - Algorithm 3 only finds one diamond and calls it a day

F1 Score - Definition

- We get these results:

Approach	Precision (%)	Recall (%)
1	60	40
2	12	98
3	95	5

- The F1 score (a.k.a the “harmonic mean” of precision and recall) provides a much more balanced measure:

$$\text{F1 Score} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} (\times 100)$$

F1 Score - Characteristics

- We get these results:

Approach	Precision (%)	Recall (%)	F1 Score (%)
1	60	40	48
2	12	98	21
3	95	5	10

- Does the F1 Score reflect what we see in precision and recall?
 - Algorithm 1 is more balanced between precision and recall
 - Algorithm 2 simply says “Dig everywhere”
 - Algorithm 3 only “finds one diamond and calls it a day”



F1 Score - Interpretation

- When precision is 1 or 100% AND recall is 1 or 100% \rightarrow F1 Score is 1 or 100%
 - In order to get a perfect F1 score, **BOTH** the precision **AND** recall have to be perfect
- If either precision is 0 OR recall is 0 \rightarrow F1 Score is 0
 - If either the precision or recall suffers (regardless of whether or not the other is high), then the F1 score suffers accordingly too
 - There's no cheating the F1 score: having a very high precision/recall value while the other precision/recall value is very low will clearly reflect in the F1 score
- F1 score favours more balanced precision and recall scores

F1 Score

- E.g. For cat detection:
 - 275 positive examples of which 200 were INcorrectly predicted
 - 260 negative examples of which 20 were INcorrectly predicted
- Precision:
 - $\frac{TP}{TP+FP} = \frac{75}{75+20} \times 100 = 78.9\%$
- Recall:
 - $\frac{TP}{TP+FN} = \frac{75}{75+200} \times 100 = 27.3\%$
- F1 Score:
 -

		predicted	
		1	0
actually	1	75	200
	0	20	240

F1 Score

- E.g. For cancer detection:
 - 160 positive examples of which 10 were INcorrectly predicted
 - 1750 negative examples of which 600 were correctly predicted
- Precision:
 - $\frac{TP}{TP+FP} = \frac{150}{150+1150} \times 100 = 11.5\%$
- Recall:
 - $\frac{TP}{TP+FN} = \frac{150}{150+10} \times 100 = 93.8\%$
- F1 Score:
 -

		predicted	
		1	0
actually	1	150	10
	0	1150	600

Model Evaluation For Classification

- In general, if we have a number of different approaches e.g. different features that we want to compare, we:
 1. Train each of the approaches on the Train set
 2. Compute the accuracy, precision, recall and F1 score on the CV set: pick the one with the highest F1 score
 3. Compute the accuracy, precision, recall and F1 score on the Test set: quote this as the final performance of the approach

THE END

Of Logistic Regression Part 1

Content

- Classification With Logistic Regression
 - Model Representation
 - Obtaining the Decision Boundary
 - Cost Function
 - Cost Function Minimization For Classification
 - Multi-Class Classification
- Model Evaluation For Classification
 - Accuracy
 - Precision and Recall
 - F1 score
- Overfitting and Underfitting
 - Introduction
 - The Concept of Regularization
 - Regularized Linear Regression
 - Regularized Logistic Regression