# DS 4400

# Machine Learning and Data Mining I

Alina Oprea

Associate Professor, CCIS

Northeastern University

January 31 2019

# Logistics

- HW 2 is due on Friday 02/08
- Tentative schedule of HW on class website
- <span style="color:red">Project proposal: due Feb 21</span>
  - 1 page description of problem you will solve, dataset, and ML algorithms
  - Individual project
  - Project template and potential ideas are on Piazza
- <span style="color:red">Project milestone: due March 21</span>
  - 2 page description on progress
- <span style="color:red">Project report at the end of semester and project presentations in class (10 minute per project)</span>
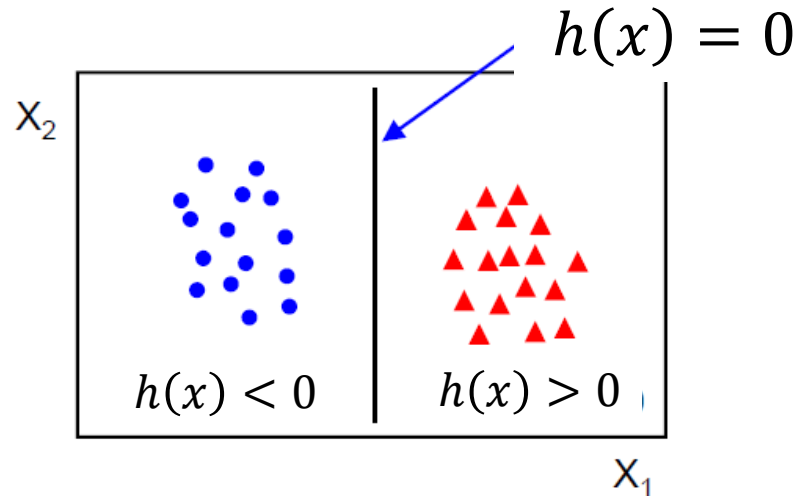
# Outline

- Logistic regression
  - Classification based on probability
- Maximum Likelihood Estimation (MLE)
  - Application to Logistic Regression
- Gradient Descent for Logistic Regression
- Evaluation of classifiers
  - Metrics
  - ROC curves

# Linear classifiers

A linear classifier has the form

$$h_\theta(x) = f(\theta^T x)$$

$$h(x) = 0$$

$X_2$

$h(x) < 0$   $h(x) > 0$

$X_1$

- Properties
  - $(\theta_0, \theta_1, \dots, \theta_d)$ = model parameters
  - Decision boundary is a hyper-plane
  - Perceptron is a special case with $f = sign$
- Pros
  - Very compact model (size d)
  - Perceptron is fast
- Cons
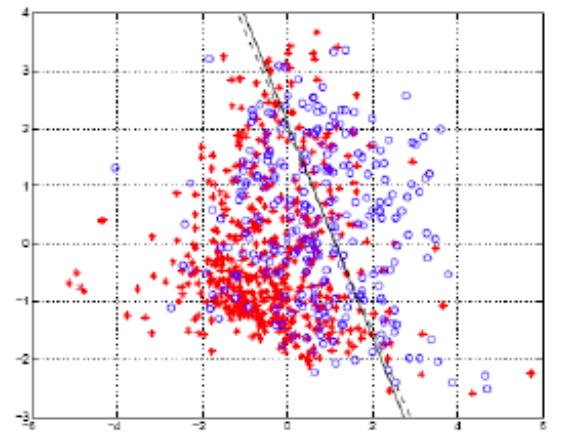  - Does not work for data that is not linearly separable

# Classification based on Probability

- Instead of just predicting the class, give the probability of the instance being that class
  - i.e., learn $p(y \mid x)$

- Comparison to perceptron:
  - Perceptron doesn't produce probability estimate

- Recall that:
$$0 \leq p(\text{event}) \leq 1$$
$$p(\text{event}) + p(\neg\text{event}) = 1$$

# Logistic regression

- Takes a probabilistic approach to learning discriminative functions (i.e., a classifier)

- $h_{\boldsymbol{\theta}}(\boldsymbol{x})$ should give $p(y = 1 \mid \boldsymbol{x}; \boldsymbol{\theta})$

  – Want $0 \leq h_{\boldsymbol{\theta}}(\boldsymbol{x}) \leq 1$

  Can't just use linear regression with a threshold

- Logistic regression model:

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = g\left(\boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{x}\right)$$

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

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

Logistic / Sigmoid Function



$g(z)$

# LR is a Linear Classifier!

- Predict $y = 1$ if:

$$\mathrm{P}[y = 1 | x; \theta] > \mathrm{P}[y = 0 | x; \theta]$$

$$\mathrm{P}[y = 1 | x; \theta] > \text{½}$$

$$\frac{1}{1 + e^{-\theta^T x}} > \frac{1}{2}$$

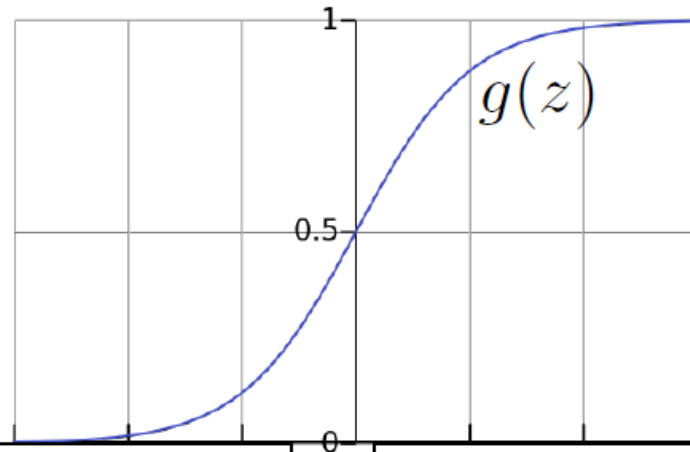- Equivalent to:

- $e^{\theta_0 + \sum_{i=1}^{d} \theta_j x_j} > 1$

- $\theta_0 + \sum_{i=1}^{d} \theta_j x_j > 0$

Logistic Regression is a linear classifier!

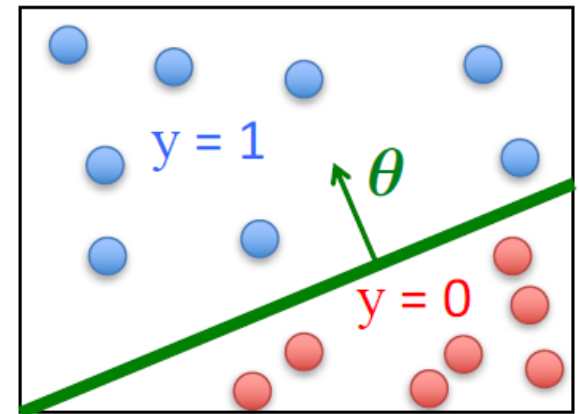# Logistic Regression

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

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

$g(z)$

$\theta^\mathsf{T} x$ should be large <u>negative</u> values for negative instances

$\theta^\mathsf{T} x$ should be large <u>positive</u> values for positive instances

- Assume a threshold and...
  - Predict y = 1 if $h_\theta(x) \geq 0.5$
  - Predict y = 0 if $h_\theta(x) < 0.5$

y = 1

$\theta$

y = 0

Logistic Regression is a linear classifier!

# Logistic Regression

- Given $\left\{ \left( \boldsymbol{x}^{(1)}, y^{(1)} \right), \left( \boldsymbol{x}^{(2)}, y^{(2)} \right), \ldots, \left( \boldsymbol{x}^{(n)}, y^{(n)} \right) \right\}$
  where $\boldsymbol{x}^{(i)} \in \mathbb{R}^d, \ y^{(i)} \in \{0, 1\}$

- Model: $h_{\boldsymbol{\theta}}(\boldsymbol{x}) = g\left( \boldsymbol{\theta}^{\mathsf{T}} \boldsymbol{x} \right)$

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

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \qquad \boldsymbol{x}^{\mathsf{T}} = \begin{bmatrix} 1 & x_1 & \ldots & x_d \end{bmatrix}$$

# Logistic Regression Objective

- Can't just use squared loss as in linear regression:

$$J(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}} \left( \boldsymbol{x}^{(i)} \right) - y^{(i)} \right)^2$$

  – Using the logistic regression model

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

  results in a non-convex optimization

# Maximum Likelihood Estimation (MLE)

Given training data $X = \{x^{(1)}, \ldots, x^{(n)}\}$ with labels $Y = \{y^{(1)}, \ldots, y^{(n)}\}$

What is the likelihood of training data for parameter $\theta$?

Define likelihood function

$$Max_{\theta} \, L(\theta) = P[Y|X; \theta]$$

Assumption: training points are independent

$$L(\theta) = \prod_{i=1}^{n} P[y^{(i)}|x^{(i)}; \theta]$$

General probabilistic method for classifier training

# Log Likelihood

- Max likelihood is equivalent to maximizing log of likelihood

$$L(\theta) = \prod_{i=1}^{n} P[y^{(i)}|x^{(i)}, \theta]$$

$$\log L(\theta) = \sum_{i=1}^{n} \log P[y^{(i)}|x^{(i)}, \theta]$$

- They both have the same maximum $\theta_{MLE}$

# MLE for Logistic Regression

$$p(y|x,\theta) = h_\theta(x)^y \left(1 - h_\theta(x)\right)^{1-y}$$

$$\boldsymbol{\theta}_{\mathrm{MLE}} = \arg\max_{\boldsymbol{\theta}} \sum_{i=1}^{n} \log p(y^{(i)} \mid \boldsymbol{x}^{(i)}; \boldsymbol{\theta})$$

$$= \arg\max_{\boldsymbol{\theta}} \sum_{i=1}^{n} y^{(i)} \log h_\theta\left(x^{(i)}\right) + (1 - y^{(i)})\log\left(1 - h_\theta(x)\right)$$

- Substitute in model, and take negative to yield

**Logistic regression objective**:

$$\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$

$$J(\boldsymbol{\theta}) = -\sum_{i=1}^{n} \left[ y^{(i)} \log h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}) + \left(1 - y^{(i)}\right) \log\left(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)})\right) \right]$$

# Objective for Logistic Regression

$$J(\boldsymbol{\theta}) = -\sum_{i=1}^{n} \left[ y^{(i)} \log h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}) + \left(1 - y^{(i)}\right) \log \left(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)})\right) \right]$$

- Cost of a single instance:

$$\text{cost}\left(h_{\boldsymbol{\theta}}(\boldsymbol{x}), y\right) = \begin{cases} -\log(h_{\boldsymbol{\theta}}(\boldsymbol{x})) & \text{if } y = 1 \\ -\log(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x})) & \text{if } y = 0 \end{cases}$$
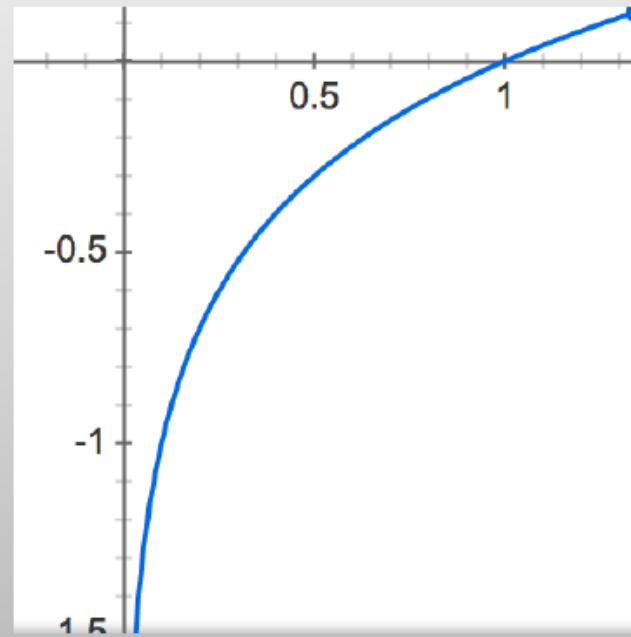
- Can re-write objective function as

$$J(\boldsymbol{\theta}) = \sum_{i=1}^{n} \text{cost}\left(h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}), y^{(i)}\right)$$

Cross-entropy loss

# Intuition

$$\text{cost}\left(h_{\boldsymbol{\theta}}(\boldsymbol{x}), y\right) = \begin{cases} -\log(h_{\boldsymbol{\theta}}(\boldsymbol{x})) & \text{if } y = 1 \\ -\log(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x})) & \text{if } y = 0 \end{cases}$$

Aside: Recall the plot of log(z)

# Intuition

$$\text{cost}\,(h_{\boldsymbol{\theta}}(\boldsymbol{x}), y) = \begin{cases} \boxed{-\log(h_{\boldsymbol{\theta}}(\boldsymbol{x})) \quad \text{if } y = 1} \\ -\log(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x})) \quad \text{if } y = 0 \end{cases}$$

If y = 1

- Cost = 0 if prediction is correct
- As $h_{\boldsymbol{\theta}}(\boldsymbol{x}) \to 0, \text{cost} \to \infty$

- Captures intuition that larger mistakes should get larger penalties
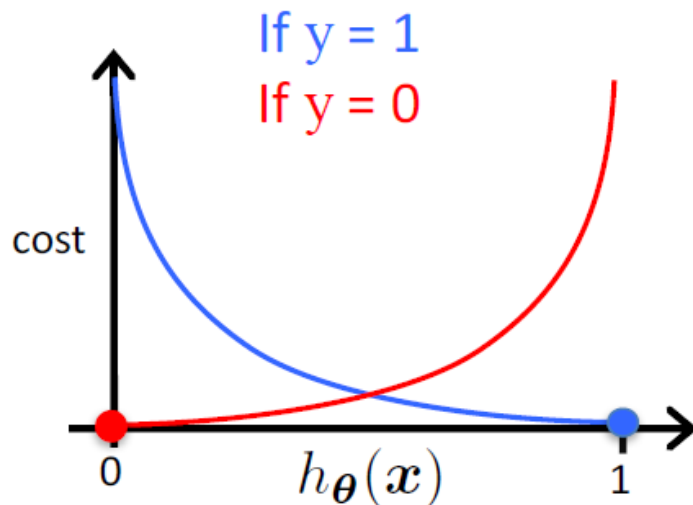  - e.g., predict $h_{\boldsymbol{\theta}}(\boldsymbol{x}) = 0$, but y = 1

If y = 1

cost

0    $h_{\boldsymbol{\theta}}(\boldsymbol{x})$    1

16

# Intuition

$$\text{cost}\left(h_{\boldsymbol{\theta}}(\boldsymbol{x}), y\right) = \begin{cases} -\log(h_{\boldsymbol{\theta}}(\boldsymbol{x})) & \text{if } y = 1 \\ -\log(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x})) & \text{if } y = 0 \end{cases}$$

If y = 0

- Cost = 0 if prediction is correct

- As $(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x})) \to 0, \text{cost} \to \infty$

- Captures intuition that larger mistakes should get larger penalties

If y = 1
If y = 0

cost

$h_{\boldsymbol{\theta}}(\boldsymbol{x})$

0          1

# Gradient Descent for Logistic Regression

$$J(\boldsymbol{\theta}) = -\sum_{i=1}^{n} \left[ y^{(i)} \log h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}) + \left(1 - y^{(i)}\right) \log \left(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)})\right) \right]$$

$$J(\theta) = -\sum_{i=1}^{n} C_i$$

Want $\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$

- Initialize $\boldsymbol{\theta}$
- Repeat until convergence

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

simultaneous update
for j = 0 ... d

# Computing Gradients

- Derivative of sigmoid
  - $g(z) = \frac{1}{1+e^{-z}}; g'(z) = \frac{e^{-z}}{(1+e^{-z})^2} = g(z)(1 - g(z))$

- Derivative of hypothesis
  - $h_\theta(x) = g(\theta^T x) = g(\theta_j x_j + \sum_{k \neq j} \theta_k x_k)$
  - $\frac{\partial h_\theta(x)}{\partial \theta_j} = \frac{\partial g(\theta^T x)}{\partial \theta_j} x_j = g(\theta^T x)(1 - g(\theta^T x))x_j$

- Derivation of $C_i$
  - $\frac{\partial C_i}{\partial \theta_j} = y^{(i)} \frac{1}{h_\theta(x^i)} g(\theta^T x^{(i)}) \left(1 - g(\theta^T x^{(i)})\right) x_j^{(i)}$ -

$$(1 - y^{(i)}) \frac{1}{1 - h_\theta(x^i)} g(\theta^T x^{(i)}) \left(1 - g(\theta^T x^{(i)})\right) x_j^{(i)}$$

$$= \left(y^{(i)} - h_\theta(x^{(i)})\right) x_j^{(i)}$$

# Gradient Descent for Logistic Regression

$$J(\boldsymbol{\theta}) = -\sum_{i=1}^{n}\left[y^{(i)}\log h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}) + \left(1 - y^{(i)}\right)\log\left(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)})\right)\right]$$

Want $\min\limits_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$

- Initialize $\boldsymbol{\theta}$

- Repeat until convergence     (simultaneous update for j = 0 ... d)

$$\theta_0 \leftarrow \theta_0 - \alpha \sum_{i=1}^{n}\left(h_{\boldsymbol{\theta}}\left(\boldsymbol{x}^{(i)}\right) - y^{(i)}\right)$$

$$\theta_j \leftarrow \theta_j - \alpha \left[\sum_{i=1}^{n}\left(h_{\boldsymbol{\theta}}\left(\boldsymbol{x}^{(i)}\right) - y^{(i)}\right)x_j^{(i)}\right]$$

# Gradient Descent for Logistic Regression

Want $\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$

- Initialize $\boldsymbol{\theta}$

- Repeat until convergence (simultaneous update for j = 0 ... d)

$$\theta_0 \leftarrow \theta_0 - \alpha \sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}}\left(\boldsymbol{x}^{(i)}\right) - y^{(i)} \right)$$

$$\theta_j \leftarrow \theta_j - \alpha \left[ \sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}}\left(\boldsymbol{x}^{(i)}\right) - y^{(i)} \right) x_j^{(i)} \right]$$

**This looks IDENTICAL to Linear Regression!**

- However, the form of the model is very different:

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

# MLE

- Probabilistic method to train classification or regression models
- Find model parameter that <span style="color:red">maximizes likelihood function</span>

$$Max_\theta \, L(\theta) = P[Y|X; \theta] = \prod_{i=1}^{n} P[y^{(i)}|x^{(i)}; \theta]$$

- Equivalent to <span style="color:red">maximize log likelihood function</span>
- Interesting property
  - MLE for linear regression has exactly the same solution as the MSE minimizer (least-square solution)

# Regularized Logistic Regression

$$J(\boldsymbol{\theta}) = -\sum_{i=1}^{n} \left[ y^{(i)} \log h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)}) + \left(1 - y^{(i)}\right) \log \left(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x}^{(i)})\right) \right]$$

- We can regularize logistic regression exactly as before:

$$J_{\text{regularized}}(\boldsymbol{\theta}) = J(\boldsymbol{\theta}) + \lambda \sum_{j=1}^{d} \theta_j^2$$

$$= J(\boldsymbol{\theta}) + \lambda \|\boldsymbol{\theta}_{[1:d]}\|_2^2$$

L2 regularization

# Classifier Evaluation

- Classification is a supervised learning problem
  - Prediction is binary or multi-class
- Classification techniques
  - Linear classifiers
    - Perceptron (online or batch mode)
    - Logistic regression (probabilistic interpretation)
  - Instance learners
    - kNN: need to store entire training data
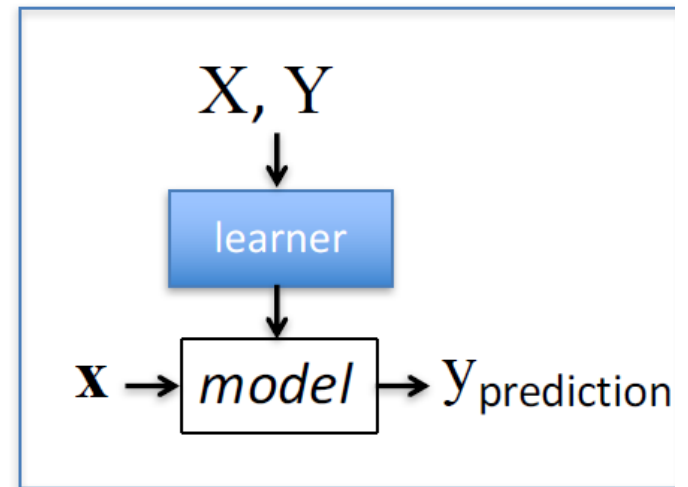- Cross-validation should be used for parameter selection and estimation of model error

# Evaluation of classifiers

**Given:** labeled training data $X, Y = \{x^{(i)}, y^{(i)}\}_{i=1}^{n}$

- Assumes each $x^{(i)} \sim \mathcal{D}(\mathcal{X})$

**Train the model:**

$model \leftarrow classifier.\text{train}(X, Y)$



**Apply the model to new data:**

- Given: new unlabeled instance $x \sim \mathcal{D}(\mathcal{X})$

$y_{\text{prediction}} \leftarrow model.\text{predict}(\mathbf{x})$

# Classification Metrics

$$\text{accuracy} = \frac{\#\ \text{correct predictions}}{\#\ \text{test instances}}$$

$$\text{error} = 1 - \text{accuracy} = \frac{\#\ \text{incorrect predictions}}{\#\ \text{test instances}}$$

- Training set accuracy and error
- Testing set accuracy and error

# Confusion Matrix

- Given a dataset of $P$ positive instances and $N$ negative instances:

**Predicted Class**

| Actual Class | | Yes | No |
|---|---|---|---|
| | Yes | TP | FN |
| | No | FP | TN |

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

- Imagine using classifier to identify positive cases (i.e., for information retrieval)

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

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

Probability that classifier predicts positive correctly

Probability that actual class is predicted correctly

# Why One Metric is Not Enough

Assume that in your training data, Spam email is 1% of data, and Ham email is 99% of data

- Scenario 1
  - Have classifier always output HAM!
  - What is the accuracy?      99%
- Scenario 2
  - Predict one SPAM email as SPAM, all other emails as legitimate
  - What is the precision?      100%
- Scenario 3
  - Output always SPAM!
  - What is the recall?      100%

# Confusion Matrix

- Given a dataset of $P$ positive instances and $N$ negative instances:

**Predicted Class**

| Actual Class | Yes | No |
|---|---|---|
| Yes | TP | FN |
| No | FP | TN |

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

- Imagine using classifier to identify positive cases (i.e., for information retrieval)

$$\text{precision} = \frac{TP}{TP + FP} \qquad \text{recall} = \frac{TP}{TP + FN}$$

$$\text{F1 score} = 2\,\frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

# Acknowledgements

- Slides made using resources from:
  - Andrew Ng
  - Eric Eaton
  - David Sontag
- Thanks!