# Day 6: Classification and Machine Learning

Kenneth Benoit

Essex Summer School 2014

July 30, 2013

# Today's Road Map

The Naive Bayes Classifier

The k-Nearest Neighbour Classifier

Support Vector Machines (SVMs)

Assessing the reliability of a training set

Evaluating classification: Precision and recall

Lab session: Classifying Text Using `quanteda`

# THE NAIVE BAYES CLASSIFIER

# Prior probabilities and updating

A test is devised to automatically flag racist news stories.

- 1% of news stories in general have racist messages
- 80% of racist news stories will be flagged by the test
- 10% of non-racist stories will also be flagged

We run the test on a new news story, and it is *flagged as* racist.

Question: What is probability that the story is *actually* racist?

Any guesses?

# Prior probabilities and updating

- What about without the test?
  - Imagine we run 1,000 news stories through the test
  - We expect that 10 will be racist
- With the test, we expect:
  - Of the 10 found to be racist, 8 should be flagged as racist
  - Of the 990 non-racist stories, 99 will be wrongly flagged as racist
  - That's a total of 107 stories flagged as racist
- So: the updated probability of a story being racist, conditional on being flagged as racist, is $\frac{8}{107} = 0.075$
- The *prior* probability of 0.01 is updated to only 0.075 by the positive test result

This is an example of Bayes' Rule:

$$P(R = 1 | T = 1) = \frac{P(T=1|R=1)P(R=1)}{P(T=1)}$$

# Multinomial Bayes model of Class given a Word

Consider $J$ word types distributed across $I$ documents, each assigned one of $K$ classes.

*At the word level*, Bayes Theorem tells us that:

$$P(c_k|w_j) = \frac{P(w_j|c_k)P(c_k)}{P(w_j)}$$

For two classes, this can be expressed as

$$= \frac{P(w_j|c_k)P(c_k)}{P(w_j|c_k)P(c_k) + P(w_j|c_{\neg k})P(c_{\neg k})} \quad (1)$$

# Classification as a goal

- Machine learning focuses on identifying classes (classification), while social science is typically interested in locating things on latent traits (scaling)
- One of the simplest and most robust classification methods is the "Naive Bayes" (NB) classifier, built on a Bayesian probability model
- The class predictions for a collection of words from NB are great for classification, but useless for scaling
- But intermediate steps from NB turn out to be excellent for scaling purposes, and identical to Laver, Benoit and Garry's "Wordscores"
- Applying lessons from machine to learning to supervised scaling, we can
  - Apply classification methods to scaling
  - improve it using lessons from machine learning

# Supervised v. unsupervised methods compared

- The goal (in text analysis) is to differentiate *documents* from one another, treating them as "bags of words"
- Different approaches:
  - *Supervised methods* require a training set that exmplify constrasting classes, identified by the researcher
  - *Unsupervised methods* scale documents based on patterns of similarity from the term-document matrix, without requiring a training step
- Relative advantage of supervised methods:
  You already know the dimension being scaled, because you set it in the training stage
- Relative disadvantage of supervised methods:
  You *must* already know the dimension being scaled, because you have to feed it good sample documents in the training stage

# Supervised v. unsupervised methods: Examples

- General examples:
  - Supervised: Naive Bayes, k-Nearest Neighbor, Support Vector Machines (SVM)
  - Unsupervised: correspondence analysis, IRT models, factor analytic approaches
- Political science applications
  - Supervised: Wordscores (LBG 2003); SVMs (Yu, Kaufman and Diermeier 2008); Naive Bayes (Evans et al 2007)
  - Unsupervised "Wordfish" (Slapin and Proksch 2008); Correspondence analysis (Schonhardt-Bailey 2008); two-dimensional IRT (Monroe and Maeda 2004)

# Multinomial Bayes model of Class given a Word

Consider $J$ word types distributed across $I$ documents, each assigned one of $K$ classes.

*At the word level*, Bayes Theorem tells us that:

$$P(c_k|w_j) = \frac{P(w_j|c_k)P(c_k)}{P(w_j)}$$

For two classes, this can be expressed as

$$= \frac{P(w_j|c_k)P(c_k)}{P(w_j|c_k)P(c_k) + P(w_j|c_{\neg k})P(c_{\neg k})} \tag{2}$$

# Multinomial Bayes model of Class given a Word
## Class-conditional word likelihoods

$$P(c_k|w_j) = \frac{P(w_j|c_k)P(c_k)}{P(w_j|c_k)P(c_k) + P(w_j|c_{\neg k})P(c_{\neg k})}$$

- The word likelihood within class
- The maximum likelihood estimate is simply the proportion of times that word $j$ occurs in class $k$, but it is more common to use Laplace smoothing by adding 1 to each oberved count within class

# Multinomial Bayes model of Class given a Word
## Word probabilities

$$P(c_k|w_j) = \frac{P(w_j|c_k)P(c_k)}{P(w_j)}$$

- This represents the word probability from the training corpus
- Usually uninteresting, since it is constant for the training data, but needed to compute posteriors on a probability scale

# Multinomial Bayes model of Class given a Word
## Class prior probabilities

$$P(c_k|w_j) = \frac{P(w_j|c_k)P(c_k)}{P(w_j|c_k)P(c_k) + P(w_j|c_{\neg k})P(c_{\neg k})}$$

- This represents the class prior probability
- Machine learning typically takes this as the document frequency in the training set
- This approach is flawed for scaling, however, since we are scaling the latent class-ness of an unknown document, not predicting class – uniform priors are more appropriate

# Multinomial Bayes model of Class given a Word
## Class posterior probabilities

$$P(c_k|w_j) = \frac{P(w_j|c_k)P(c_k)}{P(w_j|c_k)P(c_k) + P(w_j|c_{\neg k})P(c_{\neg k})}$$

- This represents the posterior probability of membership in class $k$ for word $j$
- Under *certain conditions*, this is identical to what LBG (2003) called $P_{wr}$
- Under those conditions, the LBG "wordscore" is the linear difference between $P(c_k|w_j)$ and $P(c_{\neg k}|w_j)$

# Moving to the document level

- The "Naive" Bayes model of a joint document-level class posterior assumes conditional independence, to multiply the word likelihoods from a "test" document, to produce:

$$P(c|d) = P(c) \prod_j \frac{P(w_j|c)}{P(w_j)}$$

- This is why we call it "naive": because it (wrongly) assumes:
  - *conditional independence* of word counts
  - *positional independence* of word counts

# Naive Bayes Classification Example

(From Manning, Raghavan and Schütze, *Introduction to Information Retrieval*)

▶ **Table 13.1** Data for parameter estimation examples.

|  | docID | words in document | in $c = China$? |
|---|---|---|---|
| training set | 1 | Chinese Beijing Chinese | yes |
|  | 2 | Chinese Chinese Shanghai | yes |
|  | 3 | Chinese Macao | yes |
|  | 4 | Tokyo Japan Chinese | no |
| test set | 5 | Chinese Chinese Chinese Tokyo Japan | ? |

# Naive Bayes Classification Example

**Example 13.1:** For the example in Table 13.1, the multinomial parameters we need to classify the test document are the priors $\hat{P}(c) = 3/4$ and $\hat{P}(\overline{c}) = 1/4$ and the following conditional probabilities:

$$
\begin{aligned}
\hat{P}(\text{Chinese}|c) &= (5+1)/(8+6) = 6/14 = 3/7 \\
\hat{P}(\text{Tokyo}|c) = \hat{P}(\text{Japan}|c) &= (0+1)/(8+6) = 1/14 \\
\hat{P}(\text{Chinese}|\overline{c}) &= (1+1)/(3+6) = 2/9 \\
\hat{P}(\text{Tokyo}|\overline{c}) = \hat{P}(\text{Japan}|\overline{c}) &= (1+1)/(3+6) = 2/9
\end{aligned}
$$

The denominators are $(8+6)$ and $(3+6)$ because the lengths of $text_c$ and $text_{\overline{c}}$ are 8 and 3, respectively, and because the constant $B$ in Equation (13.7) is 6 as the vocabulary consists of six terms.

We then get:

$$
\begin{aligned}
\hat{P}(c|d_5) &\propto 3/4 \cdot (3/7)^3 \cdot 1/14 \cdot 1/14 \approx 0.0003. \\
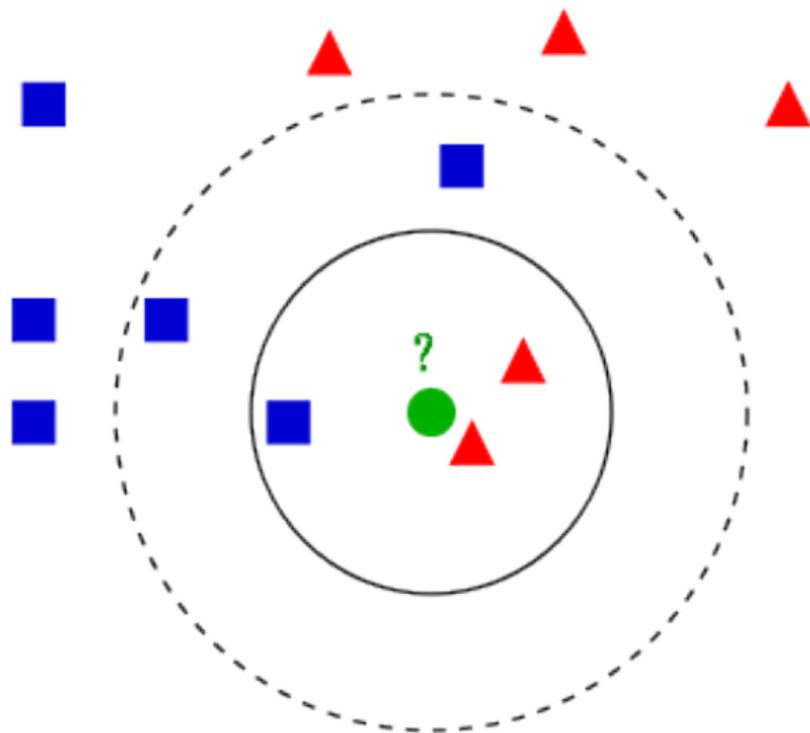\hat{P}(\overline{c}|d_5) &\propto 1/4 \cdot (2/9)^3 \cdot 2/9 \cdot 2/9 \approx 0.0001.
\end{aligned}
$$

Thus, the classifier assigns the test document to $c$ = *China*. The reason for this classification decision is that the three occurrences of the positive indicator Chinese in $d_5$ outweigh the occurrences of the two negative indicators Japan and Tokyo.
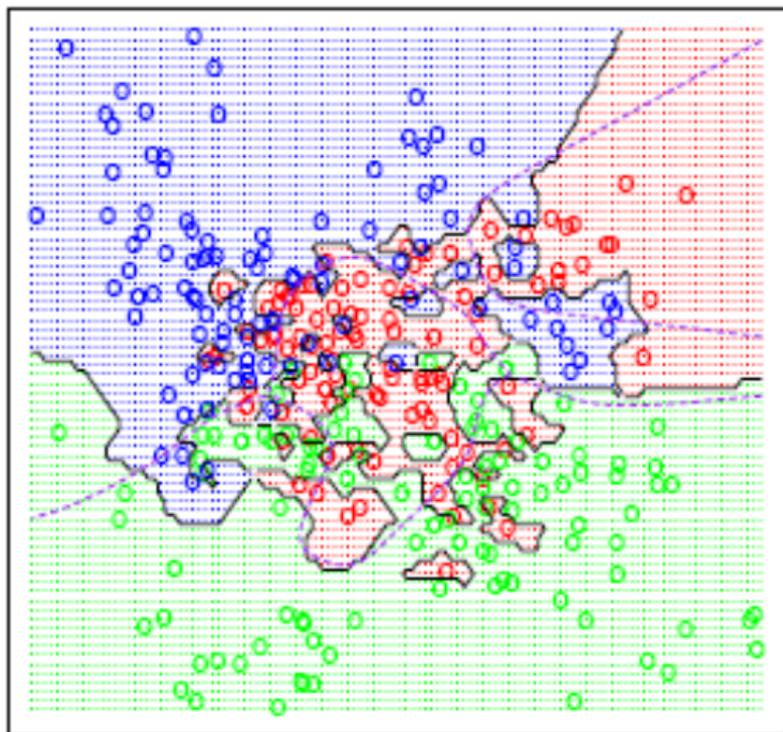
# THE *k*-NN CLASSIFIER

# Other classification methods: $k$-nearest neighbour

- A non-parametric method for classifying objects based on the training examples taht are *closest* in the feature space
- A type of instance-based learning, or "lazy learning" where the function is only approximated locally and all computation is deferred until classification
- An object is classified by a majority vote of its neighbors, with the object being assigned to the class most common amongst its $k$ nearest neighbors (where $k$ is a positive integer, usually small)
- Extremely *simple*: the only parameter that adjusts is $k$ (number of neighbors to be used) - increasing $k$ *smooths* the decision boundary

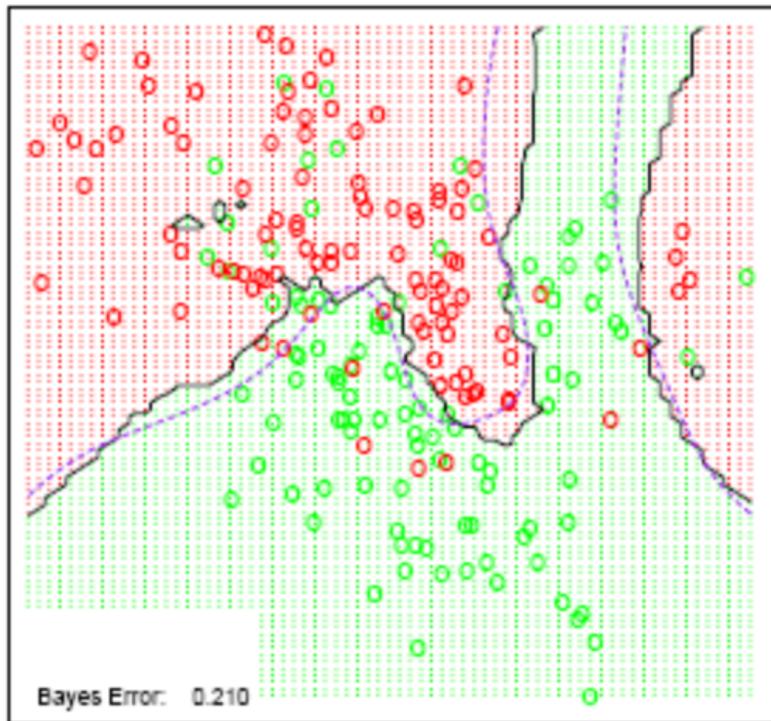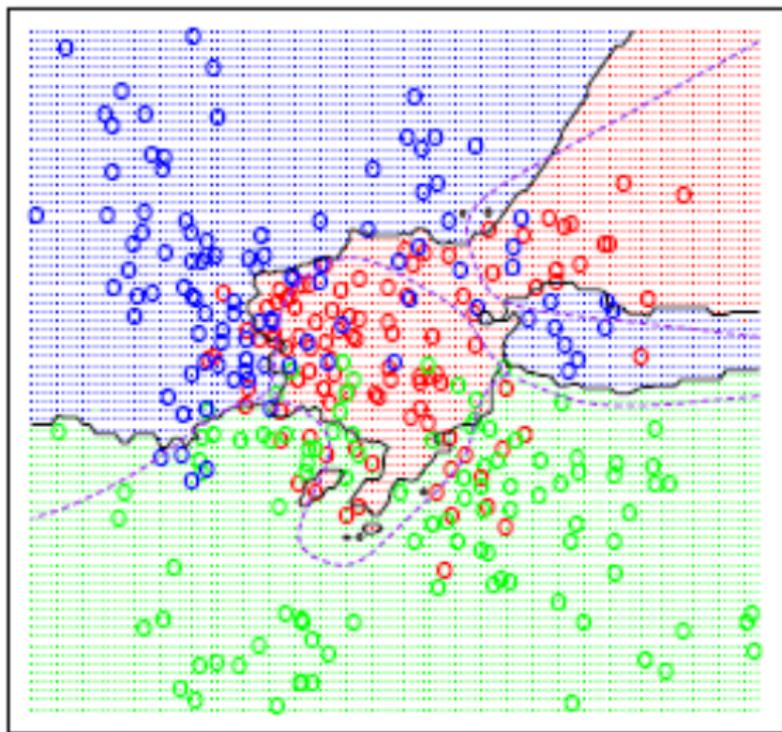# *k*-NN Example: Red or Blue?
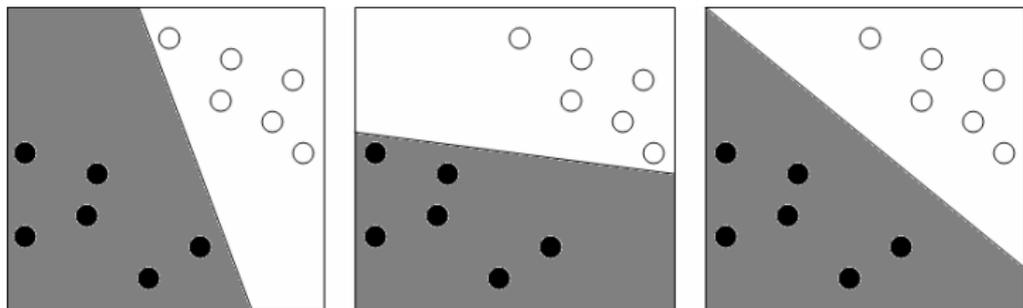
$k = 1$

$k = 7$



Bayes Error:  0.210

$k = 15$

# *k*-nearest neighbour issues: Dimensionality

- Distance usually relates to all the attributes and assumes all of them have the same effects on distance
- Misclassification may results from attributes not confirming to this assumption (sometimes called the "curse of dimensionality") – solution is to reduce the dimensions
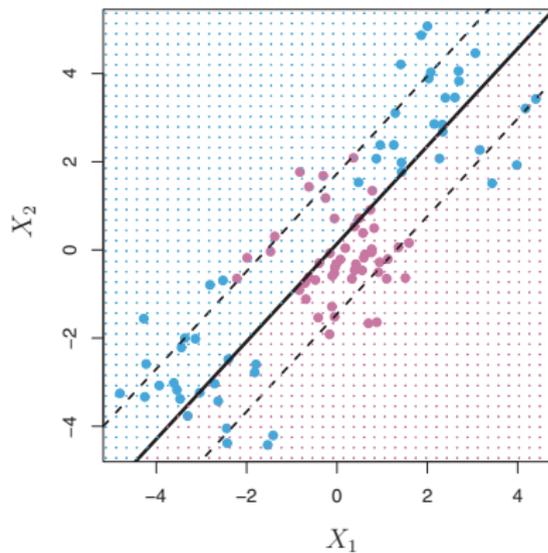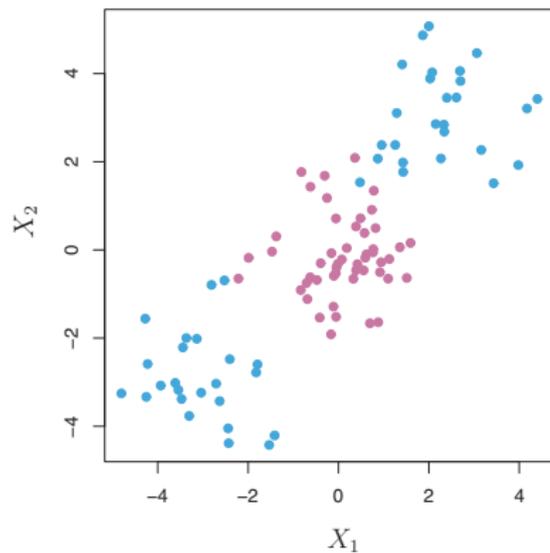- There are (many!) different *metrics* of distance

# SUPPORT VECTOR MACHINES

# (Very) General overview to SVMs

- Generalization of maximal margin classifier
- The idea is to find the classification boundary that maximizes the distance to the marginal points
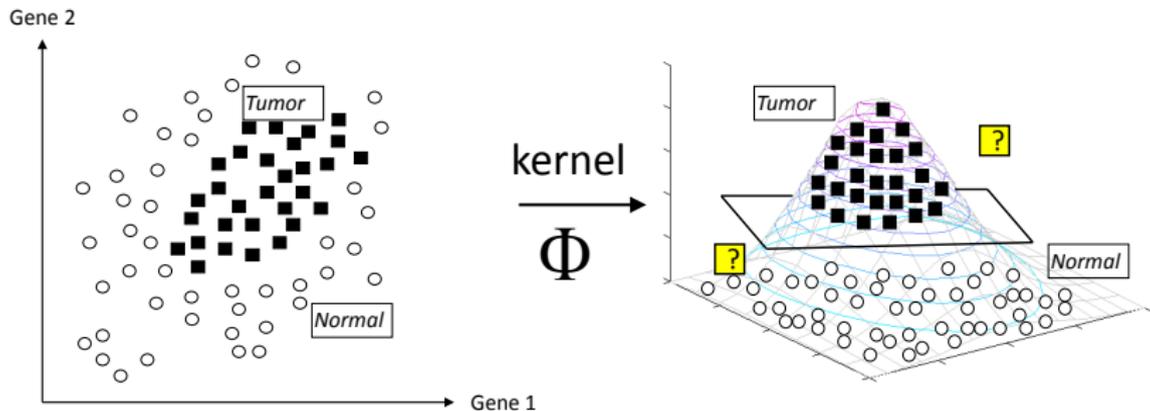


- Unfortunately MMC does not apply to cases with non-linear decision boundaries

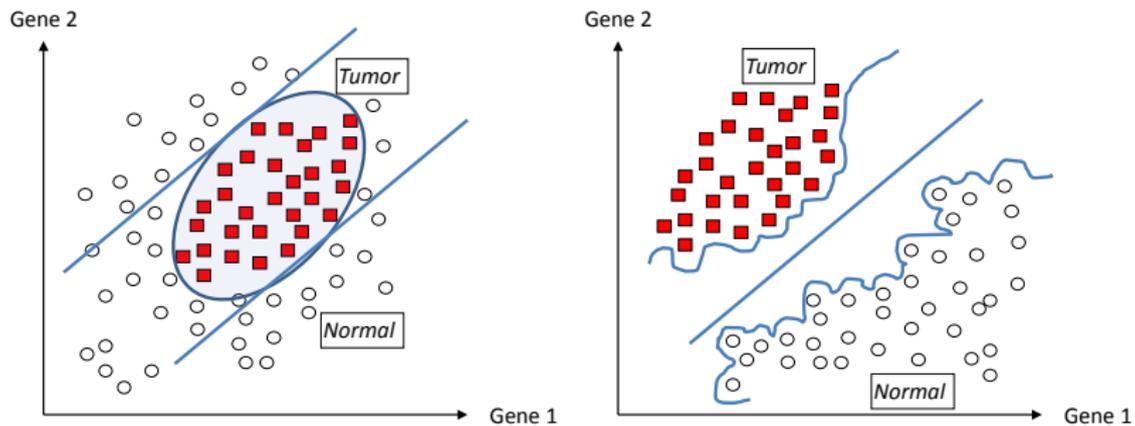# No solution to this using support vector classifier

SVMs represent the data in a *higher* dimensional projection using a kernel, and bisect this using a hyperplane



Data is not linearly separable in the underlined{input space}

Data is linearly separable in the underlined{feature space} obtained by a kernel

This is only needed when no linear separation plane exists - so not needed in second of these

# Different "kernels" can represent different decision boundaries

- ▶ This has to do with different projections of the data into higher-dimensional space
- ▶ The mathematics of this are complicated but solveable as forms of optimization problems - but the kernel choice is a user decision

# EVALUATING CLASSIFIER PERFORMANCE

# Basic principles of machine learning: Generalization and overfitting

- Generalization: A classifier or a regression algorithm learns to correctly predict output from given inputs not only in previously seen samples but also in previously unseen samples

- Overfitting: A classifier or a regression algorithm learns to correctly predict output from given inputs in previously seen samples but fails to do so in previously unseen samples. This causes poor prediction/generalization.

- Goal is to maximize the frontier of precise identification of true condition with accurate recall

# Precision and recall

- Same intuition as specificity and sensitivity earlier in course

| | | True condition | |
|---|---|---|---|
| | | Positive | Negative |
| **Prediction** | Positive | True Positive | False Positive (Type I error) |
| | Negative | False Negative (Type II error) | True Negative |

# Precision and recall and related statistics

- Precision: $\dfrac{\text{true positives}}{\text{true positives} + \text{false positives}}$

- Recall: $\dfrac{\text{true positives}}{\text{true positives} + \text{false negatives}}$

- Accuracy: $\dfrac{\text{Correctly classified}}{\text{Total number of cases}}$

- $F1 = 2\,\dfrac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$
  (the harmonic mean of precision and recall)

# Example: Computing precision/recall

Assume:

- We have a corpus where 80 documents are really positive (as opposed to negative, as in sentiment)
- Our method declares that 60 are positive
- Of the 60 declared positive, 45 are actually positive

Solution:

$$\text{Precision} = (45/(45 + 15)) = 45/60 = 0.75$$
$$\text{Recall} = (45/(45 + 35)) = 45/80 = 0.56$$

# Accuracy?

| | | True condition | |
| --- | --- | --- | --- |
| | | Positive | Negative |
| **Prediction** | Positive | 45 | | 60 |
| | Negative | | | |
| | | 80 | | |

# add in the cells we can compute



|  |  | True condition | |  |
|---|---|---|---|---|
|  |  | Positive | Negative |  |
| **Prediction** | Positive | 45 | *15* | 60 |
|  | Negative | *35* |  |  |
|  |  | 80 |  |  |

# but need True Negatives and $N$ to compute accuracy



|  |  | True condition | |
|---|---|---|---|
|  |  | Positive | Negative |
| **Prediction** | Positive | 45 | *15* | 60 |
|  | Negative | *35* | **???** |
|  |  | 80 | |

# assume 10 True Negatives:

|  |  | True condition | |  |
|---|---|---|---|---|
|  |  | Positive | Negative |  |
| **Prediction** | Positive | 45 | *15* | 60 |
|  | Negative | *35* | *10* | *45* |
|  |  | 80 | *25* | *105* |

$$\text{Accuracy} = (45 + 10)/105 \qquad\qquad = 0.52$$
$$\text{F1} = 2 * (0.75 * 0.56)/(0.75 + 0.56) \qquad = 0.64$$

## now assume 100 True Negatives:

|  |  | True condition | |  |
|---|---|---|---|---|
|  |  | Positive | Negative |  |
| **Prediction** | Positive | 45 | *15* | 60 |
|  | Negative | *35* | *100* | *135* |
|  |  | 80 | *115* | *195* |

$$\text{Accuracy} = (45 + 100)/195 \qquad\qquad = 0.74$$
$$\text{F1} = 2 * (0.75 * 0.56)/(0.75 + 0.56) \qquad = 0.64$$

# RELIABILITY TESTING FOR THE TRAINING SET

# How do we get "true" condition?

- In some domains: through more expensive or extensive tests
- In social sciences: typically by expert annotation or coding
- A scheme should be tested and reported for its reliability

# Inter-rater reliability

Different types are distinguished by the way the reliability data is obtained.

| Type | Test Design | Causes of Disagreements | Strength |
|------|-------------|-------------------------|----------|
| **Stability** | test-retest | intraobserver inconsistencies | weakest |
| **Reproducibility** | test-test | intraobserver inconsistencies + interobserver disagreements | medium |
| **Accuracy** | test-standard | intraobserver inconsistencies + interobserver disagreements + deviations from a standard | strongest |

# Measures of agreement

- **Percent agreement** Very simple: (number of agreeing ratings) / (total ratings) * 100%
- **Correlation**
  - (usually) Pearson's $r$, aka product-moment correlation
  - Formula: $r_{AB} = \frac{1}{n-1} \sum_{i=1}^{n} \left( \frac{A_i - \bar{A}}{s_A} \right) \left( \frac{B_i - \bar{B}}{s_B} \right)$
  - May also be ordinal, such as Spearman's rho or Kendall's tau-b
  - Range is [0,1]
- **Agreement measures**
  - Take into account not only observed agreement, but also *agreement that would have occured by chance*
  - Cohen's $\kappa$ is most common
  - Krippendorf's $\alpha$ is a generalization of Cohen's $\kappa$
  - Both range from [0,1]

# Reliability data matrixes

Example here used binary data (from Krippendorff)

| Article: | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|----------|---|---|---|---|---|---|---|---|---|----|
| **Coder A** | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| **Coder B** | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 |

- A and B agree on 60% of the articles: 60% agreement
- Correlation is (approximately) 0.10
- Observed *dis*agreement: 4
- Expected *dis*agreement (by chance): 4.4211
- Krippendorff's $\alpha = 1 - \frac{D_o}{D_e} = 1 - \frac{4}{4.4211} = 0.095$
- Cohen's $\kappa$ (nearly) identical