# CSE 519: Data Science
# Steven Skiena
# Stony Brook University

Lecture 22: Topics in Machine Learning

# The World of Many Weak Features

Often we have many relatively weak features to apply to a classification problem.

In text classification problems, we often have the frequency of each word in documents of positive and negative classes: e.g. the frequency of ``sale'' in spam and real email.

# Bayesian Classifiers

To classify a vector $X = (x_1, \ldots x_n)$ into one of m classes, we can use Bayes Theorem:

$$p(C_i|X) = \frac{p(C_i)p(X|C_i)}{p(X)}$$

This reduces decisions abou the class given the input to the input given the class.

# Identifying the Most Probable Class

*Argmax* is the class with the highest probability:

$$C(X) = \max_{i=1}^{m} \frac{p(C_i)p(X|C_i)}{p(X)} = \max_{i=1}^{m} p(C_i)p(X|C_i)$$

*P(Ci)* is the prior probability of class *i*.

*P(X)* is the probability of seeing input *X* over all classes.   This is dicey, but can be ignored for classification because it is constant.

# **Independence and Naive Bayes**

But what is *P(X|C)*, where *X* is a complex feature vector?

If *(a,b)* are independent, then *P(ab)=P(a) P(b)*

This calculation is much simpler than factoring in correlations and interactions of multiple factors, but:

What's the probability of having two size 9 feet?

# Complete Naive Bayes Formulation

We seek the argmax of:

$$C(X) = \max_{i=1}^{m} p(C_i)p(X|C_i) = \max_{i=1}^{m} p(C_i) \prod_{j=1}^{n} p(x_j|C_i)$$

Multiplying many probabilities is bad, so:

$$C(X) = \max_{i=1}^{m}(\log(p(C_i)) + \sum_{j=1}^{n} \log(p(x_j|C_i)))$$

# Dealing with Zero Counts

You may never have seen it before, but what is the probability my next word is defenestrate?

Observed counts do not accurately capture the frequency of rare events, for which there is typically a long tail.

Laplace asked: "What is the probability the sun will rise tomorrow?"

# +1 Discounting

Discounting is a statistical technique to adjust counts for yet-as-unseen events.

The simplest technique is add one discounting, where we add one to the frequency all outcomes, including unseen.

Thus after seeing 5 reds and 3 greens, P(new-color)=1/((5+1)+(3+1)+(0+1))

# Feature Engineering

Domain-dependent data cleaning is important:

- Z-scores and normalization
- Imputing missing values
- Dimension reduction, like SVD
- Explicit incorporation of non-linear combinations like products and ratios.

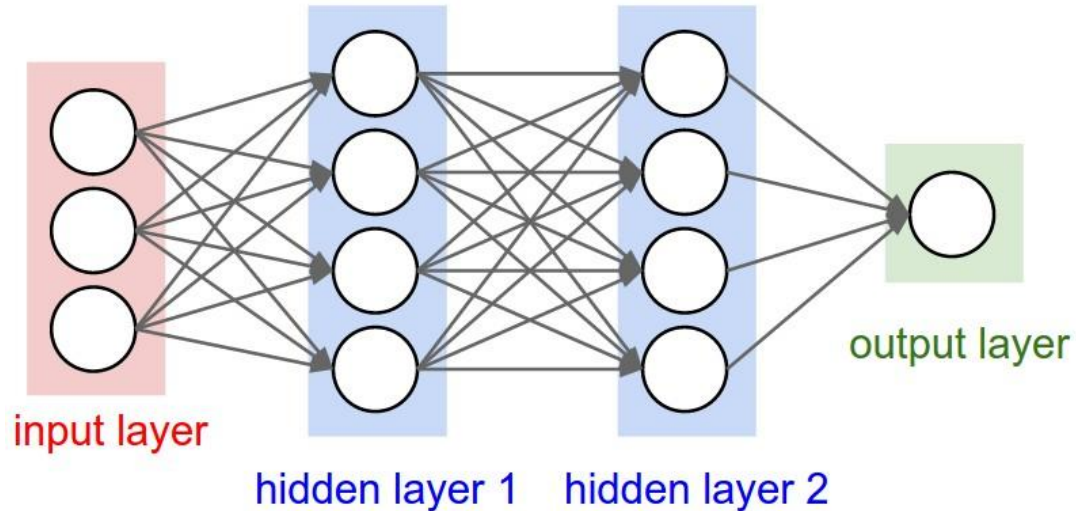# **Commissions on Art Auctions**

When you buy a painting at an auction, you pay the house a specified percentage as a fee.

How is this best represented as a feature?

- The commission percentage (e.g. *10%*)
- The actual commission paid *(0.1\*1M=$100k)*
- Change the target variable from hammer price to total amount paid: *($33M to $36.3M)*

# Deep Learning

The hottest area of machine learning today involves large, deep neural network architectures.



input layer

hidden layer 1    hidden layer 2

output layer

# **Basic Principles of Deep Learning**

- That the weight of each edge is a distinct parameter means large networks exploits large training sets.
- The depth of the networks means they can build up hierarchical representations of features: e.g. pixels, edges, regions, objects
- Toolkits like TensorFlow make it easy to build DL models if you have the data.

# Node Computations

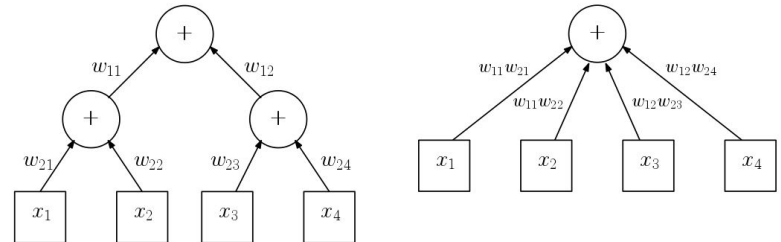Each node in the network typically computes a nonlinear function Phi(v) of a weighted input sum:
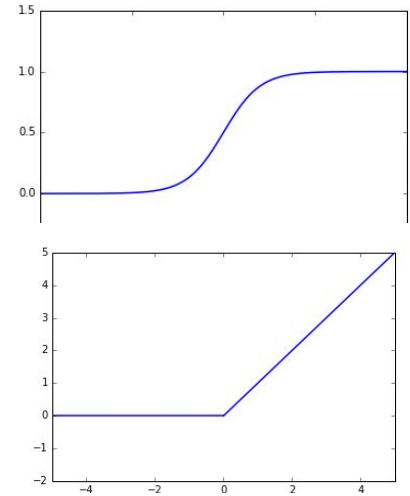
$$v_i = \beta + \sum_i w_i x_i$$

The beta term is the bias, the activation in the absence of input.

# Non-Linearity

The logit and RELU functions make good candidates for Phi.

Linear function like addition cannot exploit depth, because hidden layers add no power.

# **Backpropagation**

NNs are trained by a stochastic gradient descent-like algorithm, with changes for each training example pushed down to lower levels.

The non-linear functions result in a non-convex optimization function, but this generally produces good results.

# Word Embeddings

One NN application I have found particularly useful is word2vec, constructing 100 dimensional word representations from text corpora.

The goal is to try to predict missing words by context:     We would **** to improve

Thus large volumes of training data can be construction from text without supervision.

# Nearest Neighbors in Embeddings

| | Word | Translation |
|---|---|---|
| **French** | rouge | red |
| | juane | yellow |
| | rose | pink |
| | blanc | white |
| | orange | orange |
| | bleu | blue |

| | Word | Translation |
|---|---|---|
| **Arabic** | شكرا | thanks |
| | وشكرا | and thanks |
| | تحياتي | greetings |
| | شكراً | thanks + diacritic |
| | وشكراً | and thanks + diacritic |
| | مرحبا | hello |

| | Word | Translation |
|---|---|---|
| **Russian** | Путин | Putin |
| | Янукович | Yanukovych |
| | Троцкий | Trotsky |
| | Гитлер | Hitler |
| | Сталин | Stalin |
| | Медведев | Medvedev |

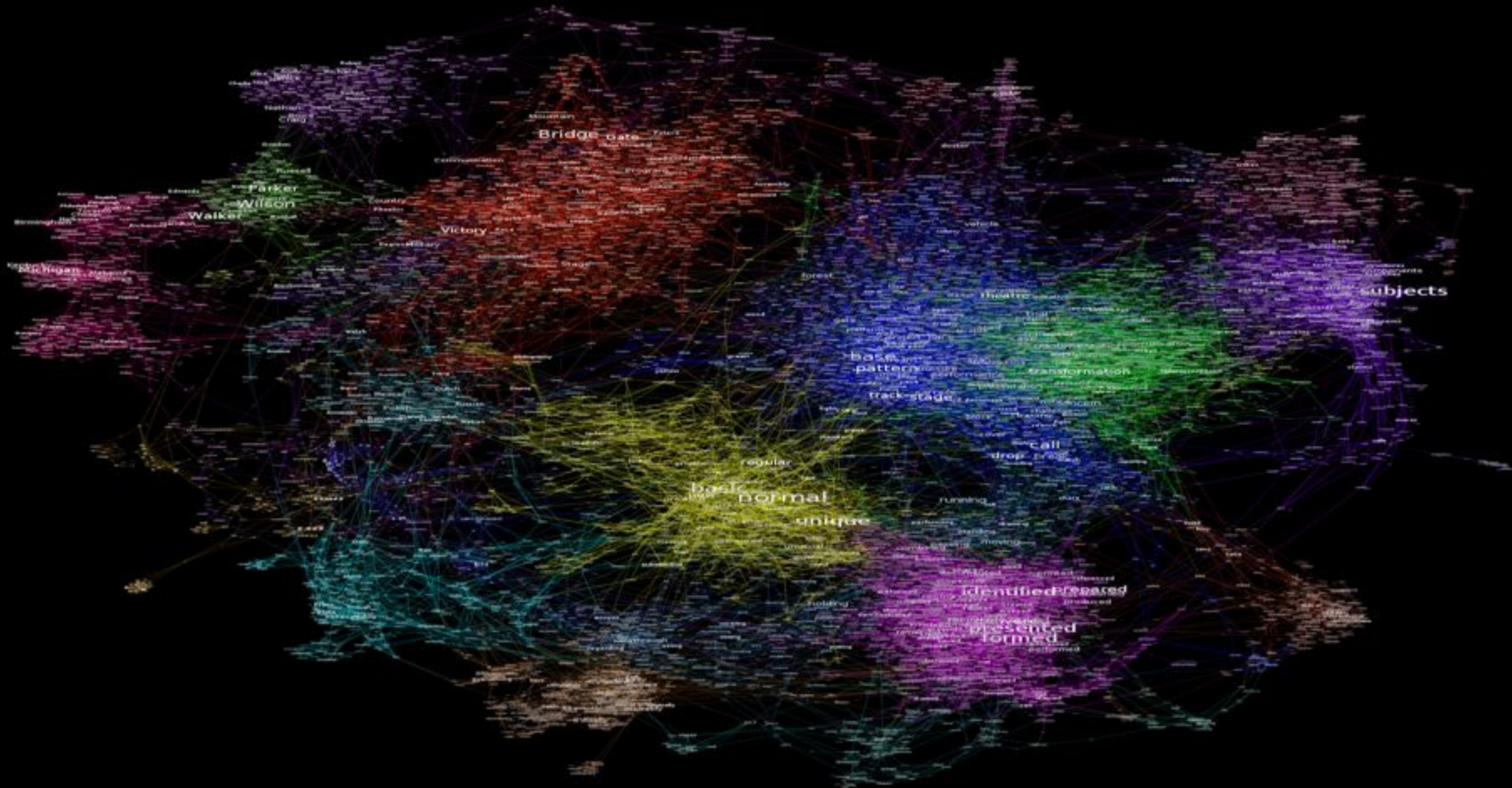| | Word | Translation |
|---|---|---|
| **Spanish** | dentista | dentist |
| | peluquero | barber |
| | ginecólog | gynecologist |
| | camionero | truck driver |
| | oftalmólogo | ophthalmologist |
| | telegrafista | telegraphist |

| | Word | Translation |
|---|---|---|
| **Arabic** | ولدان | two boys |
| | ابنان | two sons |
| | ولدين | two boys |
| | طفلان | two children |
| | ابنين | two sons |
| | ابنتان | two daughters |

| | Transliteration | |
|---|---|---|
| **Chinese** | dongzhi | Winter Solstice |
| | chunfen | Vernal Equinox |
| | xiazhi | Summer solstice |
| | qiufen | Autumnal Equinox |
| | ziye | Midnight |
| | chuxi | New Year's Eve |

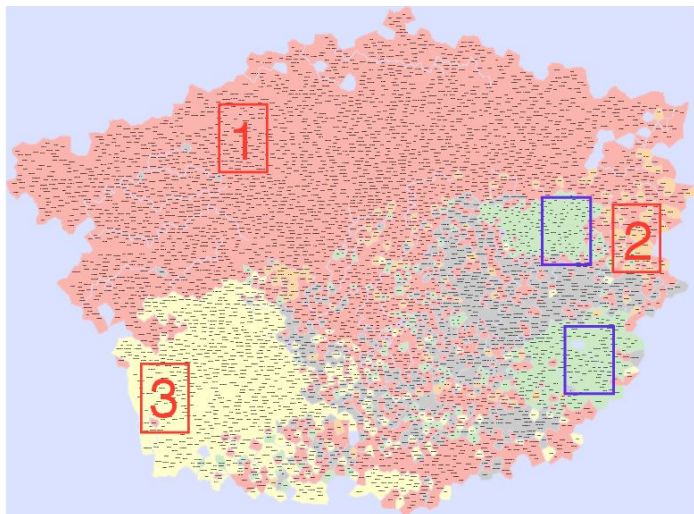| | Word | Word |
|---|---|---|
| **English** | Mumbai | Bombay |
| | Chennai | Madras |
| | Bangalore | Shanghai |
| | Kolkata | Calultta |
| | Cairo | Bangkok |
| | Hyderabad | Hyderabad |

| | Word | Word |
|---|---|---|
| **German** | Eisenbahnbetrieb | rail operations |
| | Fahrbetrieb | driving |
| | Reisezugverkehr | passenger trains |
| | Fährverkehr | ferries |
| | Handelsverkehr | Trade |
| | Schülerverkehr | students Transport |

| | Word | Word |
|---|---|---|
| **Italian** | papa | Pope |
| | Papa | Pope |
| | pontefice | pontiff |
| | basileus | basileus |
| | canridnale | cardinal |
| | frate | friar |

# Name Embeddings

Running word2vec on names from email contact lists encode gender and ethnicity:

# Graph Embeddings (DeepWalk)

Networks based on similarity or links form very sparse feature vectors.

Random walks on networks (sequences of vertices) look like sentences (sequences of words).

Thus we can use word2vec to train network representations!

# Nearest Neighbors in Wikipedia

The links between pages defines the network.

**Ludwig van Beethoven**
- Franz Schubert (0.489)
- Johannes Brahms (0.532)
- Wolfgang Mozart (0.567)
- Robert Schumann (0.576)
- Gustav Mahler (0.635)

**Mick Jagger**
- John Lennon (0.687)
- Keith Richards (0.687)
- Paul McCartney (0.796)
- Ronnie Wood (0.822)
- Eric Clapton (0.833)

**Barack Obama**
- George W. Bush (0.474)
- Hillary Clinton (0.657)
- Bill Clinton (0.658)
- Joe Biden (0.750)
- Al Gore (0.791)

**Albert Einstein**
- Richard Feynman (1.049)
- Max Planck (1.073)
- Freeman Dyson (1.107)
- Stephen Hawking (1.153)
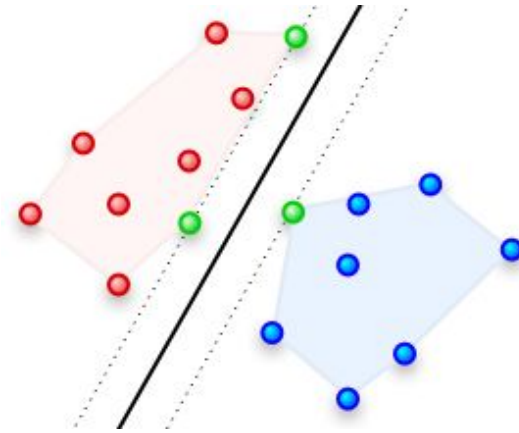- Robert Oppenheimer (1.156)

**Scarlett Johansson**
- Kirsten Dunst (0.784)
- Natalie Portman (0.786)
- Gwyneth Paltrow (0.796)
- Brad Pitt (0.858)
- Cameron Diaz (0.891)

**Steven Skiena**
- Larry Page (1.597)
- Sergey Brin (1.598)
- Danny Hillis (1.644)
- Andrei Broder (1.652)
- Mark Weiser (1.653)

# **Support Vector Machines**

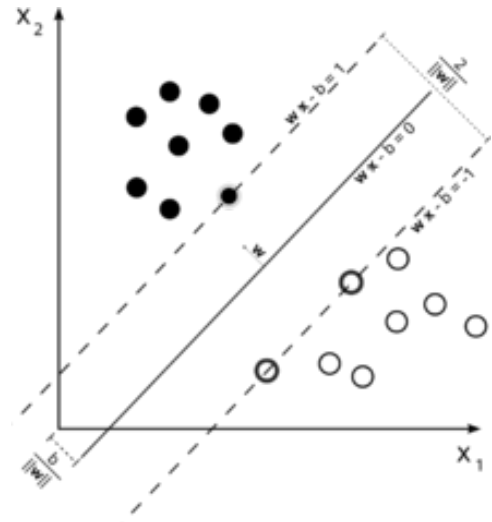SVMs are an important way to build non-linear classifiers.



They work by seeking maximum margin linear separators between the two classes.

# Optimization Problem

Optimize the coefficient size $\|\mathbf{w}\|$ subject to the constraints $y_i(\mathbf{w} \cdot \mathbf{x_i} - b) \geq 1$ for all $i = 1, \ldots, n$
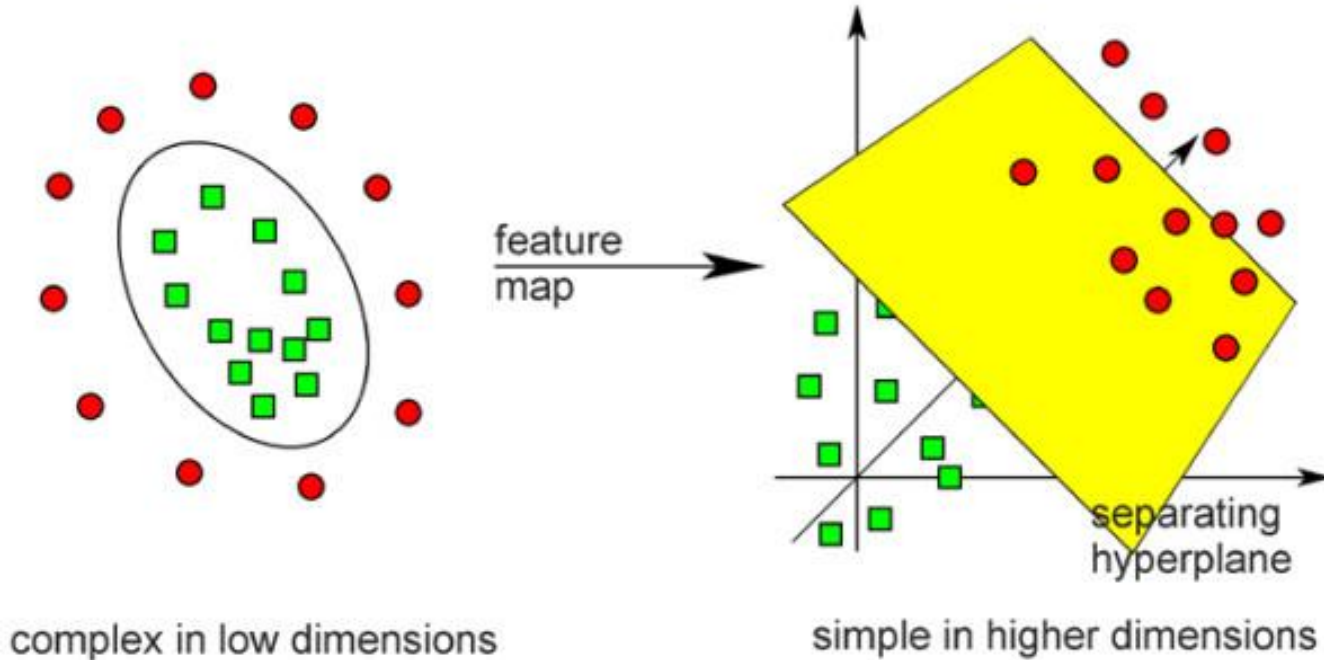
Note that only a few points (the support vectors) touch the boundary of the separating channel.

Efficient solvers like LibSVM are available for this.

# Projecting to Higher Dimensions



Separation may be easier in higher dimensions

feature map

separating hyperplane

complex in low dimensions

simple in higher dimensions

# Projecting to Higher Dimensions

The non-linearity depends upon how the space is projected to higher dimensions.

We can use features the distance from each of the n input points to the target to create an n-dimensional feature vector.

# Kernals

The magic of SVMs is that this distance matrix need not be computed explicitly.

Further, certain functions (or kernals) can be computed efficiently on these points, thus changing the feature set to yield more relevant separators.