# Machine Learning Basics

Neal Parikh

Learn 2 Quant Conference, New York
November 16, 2018

# Finding spy planes

## US Federal Agents Flew A Secret Spy Plane To Hunt Drug Cartel Leaders In Mexico

Neither the US Marshals Service nor the Mexican government wants to talk about their joint efforts to hunt drug kingpins. But BuzzFeed News spotted a Marshals spy plane circling around the time of a prominent capture in Sinaloa.

Posted on August 3, 2017, at 8:00 a.m.

**Peter Aldhous**
BuzzFeed News Reporter

**Karla Zabludovsky**
BuzzFeed News Reporter

in August 2017, Buzzfeed News publishes articles finding

- military contractors flying over SF Bay Area

- secret US Marshals plane hunting drug cartel kingpins in Mexico

- Air Force special operations planes flying over US

- . . .

# Finding spy planes



**BuzzFeed News Trained A Computer To Search For Hidden Spy Planes. This Is What We Found.**

From planes tracking drug traffickers to those testing new spying technology, US airspace is buzzing with surveillance aircraft operated for law enforcement and the military.

1. pull a publicly available dataset (not intended for this purpose)
2. train a simple machine learning model
3. validate (here, 'do journalism')

# Finding spy planes



**BuzzFeed News Trained A Computer To Search For Hidden Spy Planes. This Is What We Found.**

From planes tracking drug traffickers to those testing new spying technology, US airspace is buzzing with surveillance aircraft operated for law enforcement and the military.

1. pull 4 months of flight-tracking data from website Flightradar24
2. extract 'features': turning rates, speeds, altitudes, manufacturers
3. train a binary classifier to distinguish between previously identified FBI/DHS planes and not
4. validate

# Examples

- Adobe (font recognition using phone camera)
- Amazon (speculative shipping, Kindle browser prefetching)
- American Express (fraud detection, individual credit limits)
- Cheesecake Factory (predict food ingredient demand)
- C-SPAN (automatically name politicians on screen)
- HireVue (video analysis of job interviews for hiring/screening)
- Nest Thermostat (embedded control of smart thermostat)
- Target (market research, individualized product catalogues)
- USPS (handwriting recognition)
- Walmart (inventory, product placement)

# Automated sepsis detection for hospital operations

- sepsis is #3 leading cause of death in US, but hospitals often miss early signs and don't catch it until it's too late

- university hospitals (Duke, Johns Hopkins) deploying ML systems, some this month (Sepsis Watch), for automated sepsis detection

- *e.g.*, Duke system trained on 50K patient records, over 32M data points, with many variables (vital signs, lab tests, medical history)

- pulls patient data every 5 min to evaluate conditions, then alerts nurses

- nurses make decisions about alert, and if approved, are guided through checklist of actions

# What is machine learning?

- no precise technical definition

- usage evolved over time

- 'classical' usage is as a sub-discipline of AI research

# What is machine learning?

- intersection of computer science and statistics

- computationally tractable algorithms that learn from data

- the mathematical foundation of modern AI, but now also used in a huge variety of other domains

# What is machine learning?

- modern usage: how to build *learning procedures*, *i.e.*, how to use historical data to build a *prediction rule*

- prediction rule: algorithm mapping observable inputs to prediction of unknown quantity (the *response*)

- focus is on making predictions, and doing well on data you *haven't yet seen* (how to select the right prediction rule among several)

# What is machine learning?

- modern usage: how to build *learning procedures*, *i.e.*, how to use historical data to build a *prediction rule*

- prediction rule: algorithm mapping observable inputs to prediction of unknown quantity (the *response*)

- focus is on making predictions, and doing well on data you *haven't yet seen* (how to select the right prediction rule among several)

- informally, is mostly interchangeable with the terms 'AI' and 'modern statistical prediction' (*e.g.*, Sepsis Watch can be called 'an AI')

# Machine learning and AI

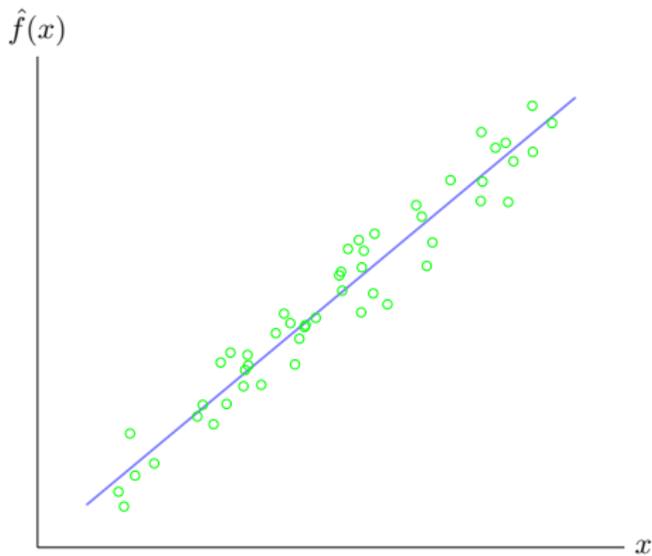| | |
|---|---|
| 1950s | Dartmouth conferences; chess & checkers; LISP; perceptron |
| 1960s | early foundational & philosophical work; formal logic |
| 1970s | neural networks; AI winter |
| 1980s | expert systems; AI winter |
| 1990s | probabilistic revolution; graphical models; kernel methods |
| 2000s | convex optimization; continuing development from 90s |
| 2010s | deep learning; large-scale & widespread applications |

## Machine learning and statistics

(Wasserman; Tibshirani)

| statistics | computer science |
| --- | --- |
| estimation/fitting | learning |
| regression/classification | supervised learning |
| clustering/density estimation | unsupervised learning |
| data | training sample |
| covariates | features, inputs |
| response | outputs |
| test set performance | generalization ability |

# Linear regression
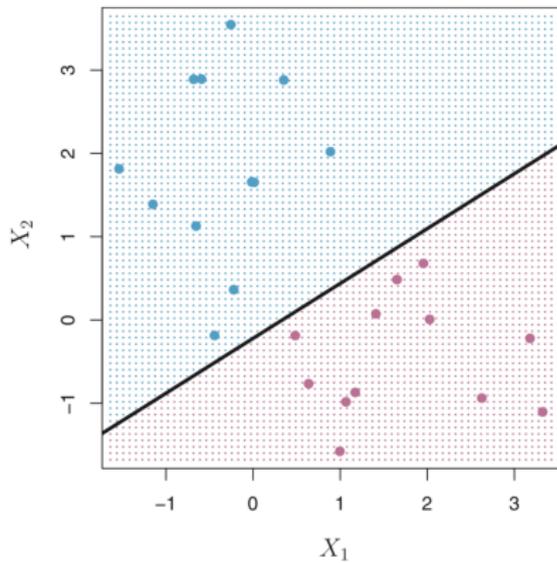


Straight line fit to 50 points in a plane.
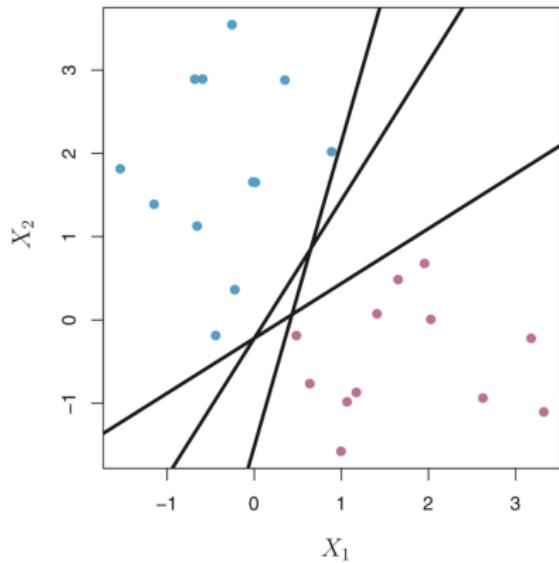
# Autoregressive time series



Hourly temperature at LAX.
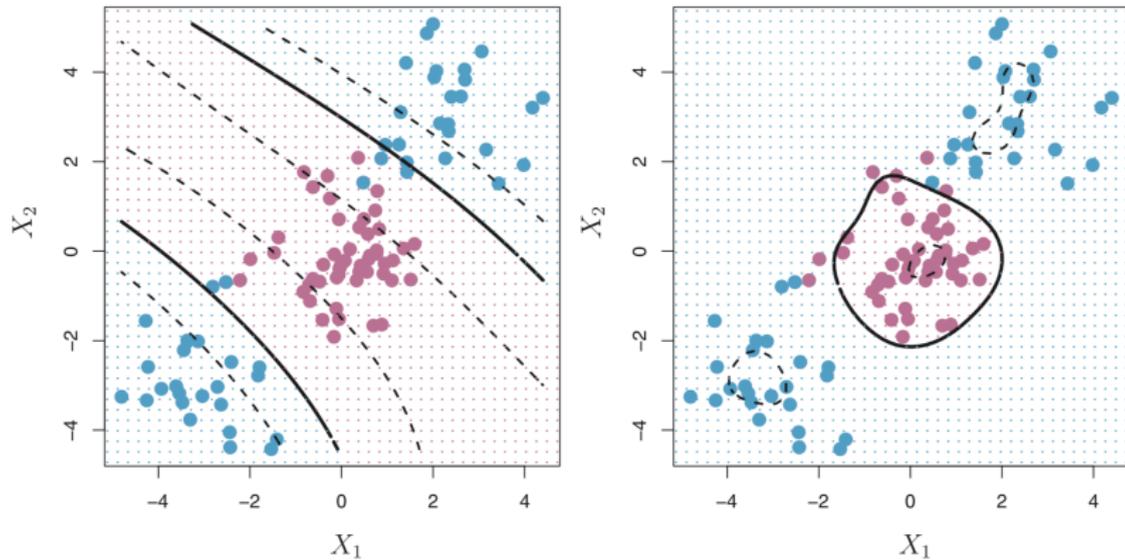
# Polynomial regression



Least squares fits of degree 2, 6, 10, and 15 to 100 points.
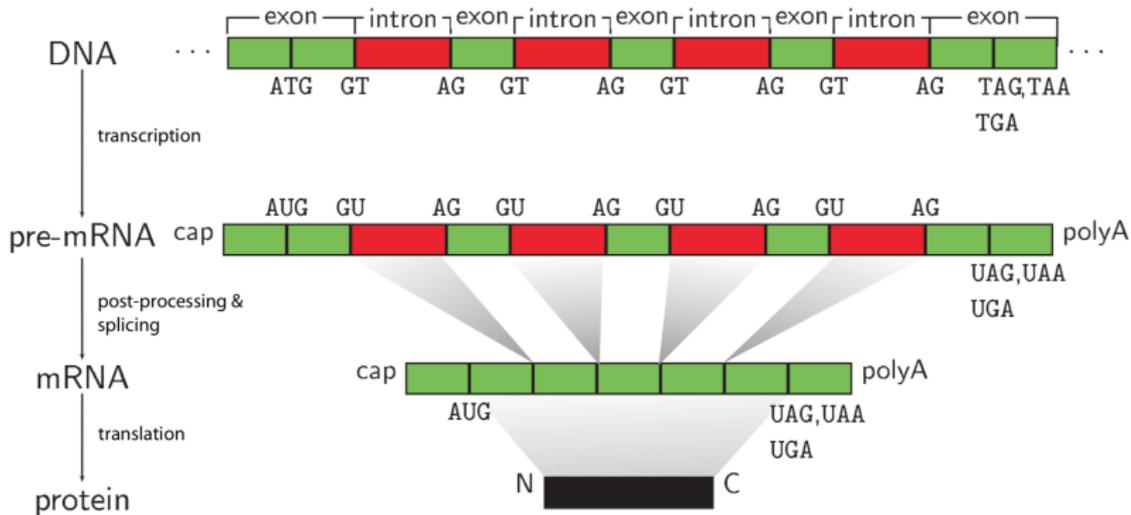
# Support vector machine
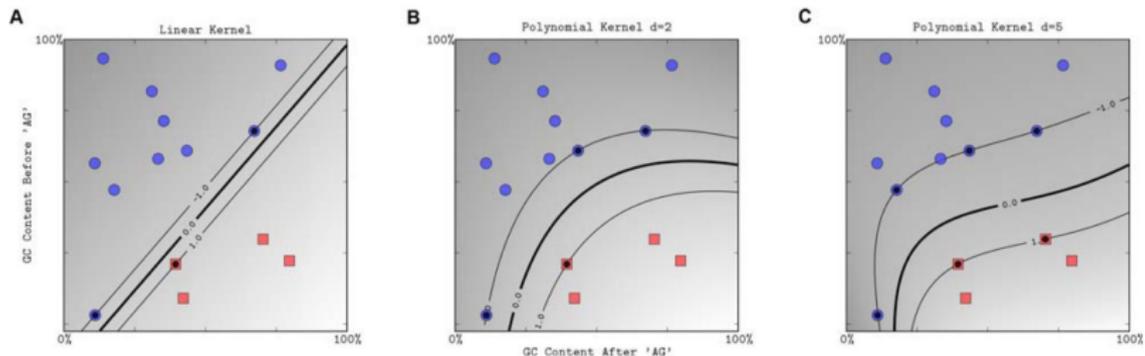
# Support vector machine

# Splice site recognition

(Ben-Hur et al., *PLoS Computational Biology*, 2008)

# Splice site recognition

(Ben-Hur et al., *PLoS Computational Biology*, 2008)

# Uses and pitfalls

uses:

- explore new, richer, unused datasets (text, image, . . . )

- internal operations (anomaly detection, data processing, . . . )

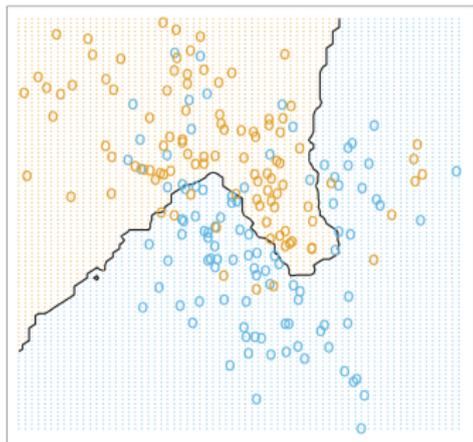- actual trading signals, portfolio construction, . . .

pitfalls:

- need appropriate team and workflows (*e.g.*, model diagnostics)

- 'bias' and ethics
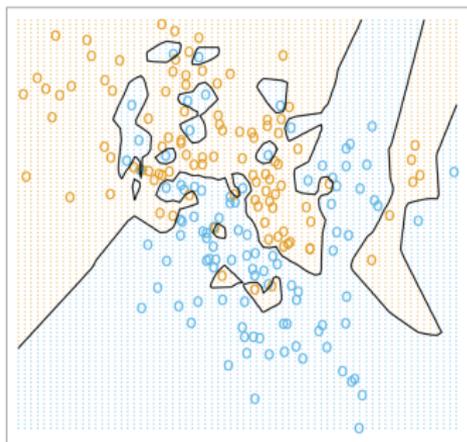
- (wrongly) anthropomorphizing models

# Thanks

Questions?

# $k$-nearest neighbors



$k = 15$        $k = 1$

Source: Hastie, Tibshirani, Friedman, *The Elements of Statistical Learning*

# Topic models

## Topic models

| sound | quantum | brain | computer | ice |
|---|---|---|---|---|
| speech | laser | memory | data | climate |
| acoustic | light | human | information | ocean |
| language | optical | visual | problem | sea |
| sounds | electron | cognitive | computers | temperature |
| stars | research | materials | fossil | volcanic |
| universe | national | organic | species | years |
| galaxies | science | molecules | evolution | fig |
| astronomers | new | molecular | birds | deposits |
| star | funding | polymer | evolutionary | rocks |

# Topic models