# Learning to Generalize

Jake Hofman

Microsoft Research

March 16, 2018

# Learning by example

# Learning by example



- How did you solve this problem?
- Can you make this process explicit (e.g. write code to do so)?

# Machine Learning: A Definition

The study of algorithms that enable machines to learn from experience and improve over time

`http://www.cs.cmu.edu/~tom/`

# Everything old is new again[1]

- Many fields ...
    - Statistics
    - Pattern recognition
    - Data mining
    - Machine learning
- ... similar goals
    - Extract and recognize patterns in data
    - Interpret or explain observations
    - Test validity of hypotheses
    - Efficiently search the space of hypotheses
    - Design efficient algorithms enabling machines to learn from data

---

[1]http://bit.ly/mloldnew

# Statistics vs. machine learning[2]

**Glossary**

| Machine learning | Statistics |
| --- | --- |
| network, graphs | model |
| weights | parameters |
| learning | fitting |
| generalization | test set performance |
| supervised learning | regression/classification |
| unsupervised learning | density estimation, clustering |
| large grant = $1,000,000 | large grant= $50,000 |
| nice place to have a meeting: Snowbird, Utah, French Alps | nice place to have a meeting: Las Vegas in August |

# Example results: machine learning[3]

Add an asymmetric frequency feature $\mathbf{y}_{j,f_{ut}}^{(3)}$: **SBRAMF-UTB-UTF-MTF-ATF-MFF-AFF**

$$\widehat{r_{uit}} = \mu_i + \mu_u + \mu_{u,t} + \mu_{i,\mathrm{bin}(t)} + \left(\mathbf{p}_i^{(1)} + \mathbf{p}_{i,\mathrm{bin}(t)}^{(2)} + \mathbf{p}_{i,f_{ut}}^{(3)}\right)^T \left(\mathbf{q}_u^{(1)} + \mathbf{q}_{u,t}^{(2)} + \frac{1}{\sqrt{|N(u)|}} \sum_{j \in N(u)} \left(\mathbf{y}_j^{(1)} + \mathbf{y}_{j,\mathrm{bin}(t)}^{(2)} + \mathbf{y}_{j,f_{ut}}^{(3)}\right)\right)$$

$$(34)$$

| Model extension (+) | epoch time | #epochs | probeRMSE, $k = 50$ features |
|---|---|---|---|
| **SBRMF** - SVD with biases | 17[s] | 69 | 0.9054 |
| **SBRAMF** - asymmetric part | 50[s] | 30 | 0.8974 |
| **+UTB** - user time bias | 61[s] | 50 | 0.8919 |
| **+UTF** - user time feature | 62[s] | 38 | 0.8911 |
| **+MTF** - movie time feature | 74[s] | 37 | 0.8908 |
| **+ATF** - asymmetric time feature | 74[s] | 44 | 0.8905 |
| **+MFF** - movie frequency feature | 149[s] | 46 | 0.8900 |
| **+AFF** - asymmetric frequency feature | 206[s] | 45 | 0.8886 (0.8846 with $k = 1000$) |

---

[3]Bell, Koren & Volinsky, 2008

# Example results: social science[4]

Table 2: Relationship of clicks/added revisions and time dummies for direct neighbors of shocked articles in the 'featured articles' condition.

| | clicks | | | $\Delta$ revisions | | |
|---|---|---|---|---|---|---|
| | (1) compare control | (2) comp. placebo | (3) before after | (4) comp. control | (5) comp. placebo | (6) before after |
| Before: days - 7 to -4 | 0.821 (1.342) | 0.486 (1.222) | -0.181 (1.066) | -0.000 (0.004) | 0.002 (0.004) | 0.001 (0.003) |
| Before: days - 3 to -1 | 1.811 (2.408) | 2.026 (1.990) | 1.910 (1.462) | 0.000 (0.005) | -0.001 (0.004) | -0.001 (0.003) |
| **t = 0** | 28.546*** (5.948) | 34.577*** (5.731) | 31.603*** (5.573) | 0.030*** (0.009) | 0.033*** (0.008) | 0.030*** (0.007) |
| t = 1 | 1.632 (2.146) | 1.535 (2.318) | 0.974 (1.565) | 0.007 (0.007) | 0.006 (0.007) | 0.005 (0.005) |
| t = 2 | -0.569 (2.768) | -1.189 (2.395) | 0.028 (1.910) | -0.013* (0.007) | -0.011 (0.007) | -0.007* (0.004) |
| After: days 3 to 6 | -2.170 (2.052) | -0.376 (2.296) | -0.531 (1.359) | 0.002 (0.004) | -0.000 (0.004) | -0.001 (0.003) |
| After: days 7 to 14 | -0.639 (2.593) | 0.207 (2.794) | 0.207 (1.953) | 0.001 (0.007) | -0.001 (0.007) | 0.001 (0.005) |
| Time Dummies | Yes | Yes | No | Yes | Yes | No |
| Mean dep. Variable | 36.208 | 36.559 | 37.276 | 0.045 | 0.045 | 0.047 |
| Observations | 346104 | 371382 | 186384 | 346104 | 371382 | 186384 |
| Number of Pages | 15732 | 16881 | 8472 | 15732 | 16881 | 8472 |
| Adj. $R^2$ | 0.002 | 0.003 | 0.004 | 0.000 | 0.000 | 0.000 |

---

[4]Krummer, 2015

# Social science vs. machine learning?[5]

Machine Learning

$$\hat{y}$$

Predict

~~vs~~
**and**

Social science

$$\hat{\beta}$$

Explain

Important to view prediction and explanation as compliments,
not substitutes

[5]Mullainathan & Spiess, JEP 2017

# Why Machine Learning?[6]

- Prediction is often useful in its own right
  e.g., forecasting economic outcomes, quantifying risk
- Prediction can help fill in missing data
  e.g., inferring online demographics
- Prediction can aid in causal inference
  e.g., matching, instrumental variables, heterogenous effects
- Predictive models can provide benchmarks for causal theories
  e.g., gap in performance indicates work to be done

# Outline for the day

What we'll cover:

### Concepts:
- Generalization error
- Overfitting
- Cross-validation
- Regularization

### Methods:
- k-nearest neighbors
- Naive Bayes
- Ridge regression
- Gradient descent

What we won't cover:
tree-based methods, deep neural nets, ensemble methods, unsupervised learning, reinforcement learning, . . .

# Learning by example



- How did you solve this problem?
- Can you make this process explicit (e.g. write code to do so)?

# Roadmap?

Step 1: Have data

Step 2: ???

Step 3: Profit

# Roadmap, take two

1. Get data

# Roadmap, take two

1. Get data
2. Visualize/perform sanity checks
3. Clean/filter observations
4. Choose features to represent data

# Roadmap, take two

1. Get data
2. Visualize/perform sanity checks
3. Clean/filter observations
4. Choose features to represent data
5. Specify model
6. Specify loss function

# Roadmap, take two



1. Get data
2. Visualize/perform sanity checks
3. Clean/filter observations
4. Choose features to represent data
5. Specify model
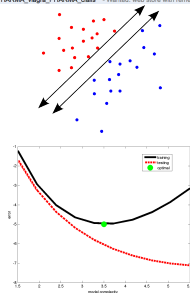6. Specify loss function
7. Develop algorithm to minimize loss

# Roadmap, take two

1. Get data
2. Visualize/perform sanity checks
3. Clean/filter observations
4. Choose features to represent data
5. Specify model
6. Specify loss function
7. Develop algorithm to minimize loss
8. Choose performance measure
9. "Train" to minimize loss
10. "Test" to evaluate generalization

# Themes

Fitting models is a relatively small piece of the pipeline

# Themes

Cleaning and normalizing data is a substantial amount of the work
(and likely impacts results)

# Themes

The features you choose to represent the data often matter more
than the algorithm that learns from it

# Themes

There's a skill in matching tools (e.g., models and algorithms) to problems

# Themes

Our models should be complex enough to explain the past, but
simple enough to generalize to the future

# Bigger models $\neq$ Better models

# Themes

Simple methods (e.g., linear models) work surprisingly well,
especially with lots of data

# Themes

Even with simple methods, there are many decisions to be made in developing a successful machine learning solution

# Digit recognition

Classification is an *supervised* learning task by which we aim to *predict the correct label* for an example given its features
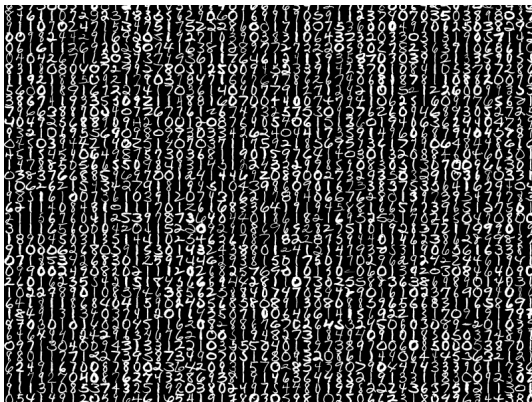


↓

0  5  4  1  4  9

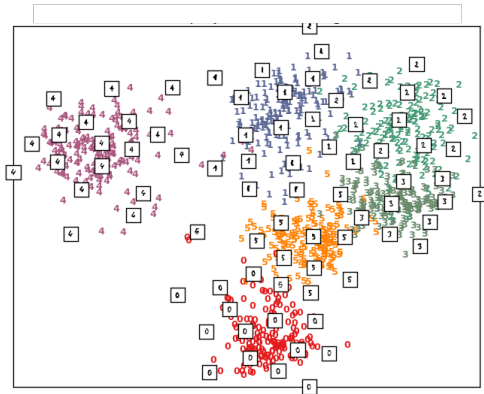e.g. determine which digit $\{0, 1, \ldots, 9\}$ is in depicted in each image

# Digit recognition

Determine which digit $\{0, 1, \ldots, 9\}$ is in depicted in each image

# k-Nearest Neighbors classification

**Memorize** training examples, predict labels using labels of the **k closest** training points



Intuition: **nearby** points have **similar** labels

# k-Nearest Neighbors classification

```
Training:

  Load all training examples into memory

Prediction:

  Compute the distance between each training
  example and the query point

  Find the k closest training examples
  to the query point

  Return a majority vote over their labels
```
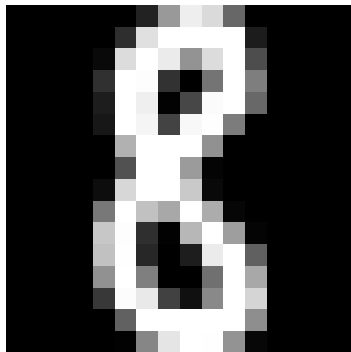
# Choosing features: Images as arrays

Grayscale images $\leftrightarrow$ 2-d arrays of $M$-by-$N$ pixel intensities
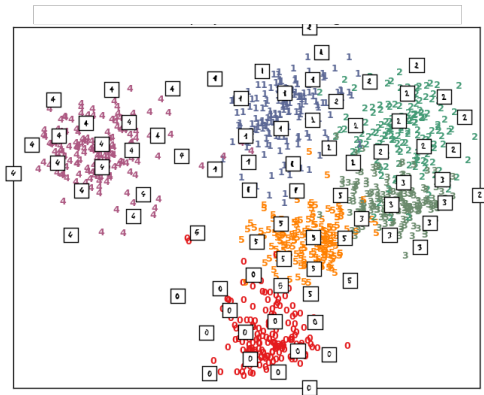
# Choosing features: Images as arrays

Grayscale images $\leftrightarrow$ 2-d arrays of $M$-by-$N$ pixel intensities



Flatten each array into a vector, representing each image as a "vector of pixels"
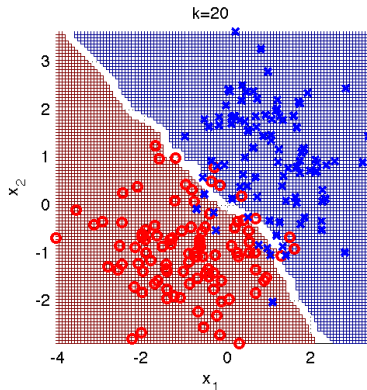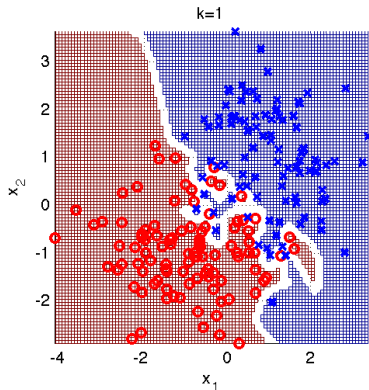
# Measuring similarity: Euclidean distance

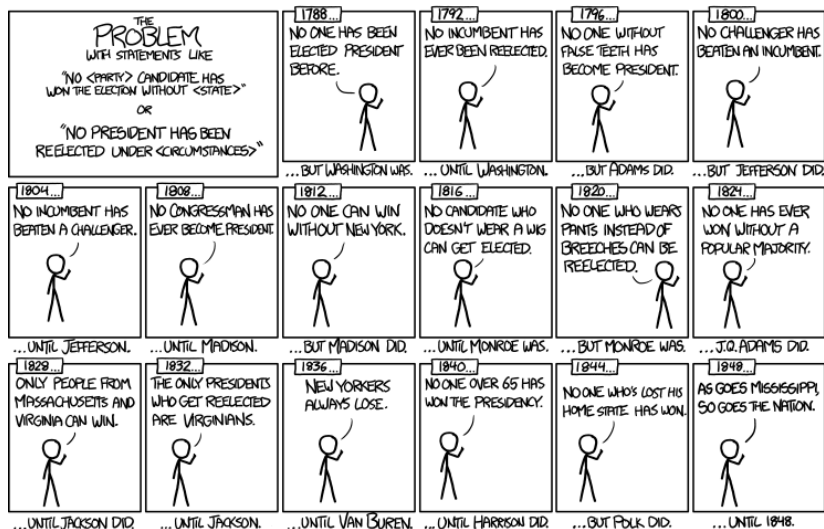Grayscale images $\leftrightarrow$ 1-d vector of $M \cdot N$ pixel intensities



Measure similarity using the Euclidean distance between any two pixel vectors
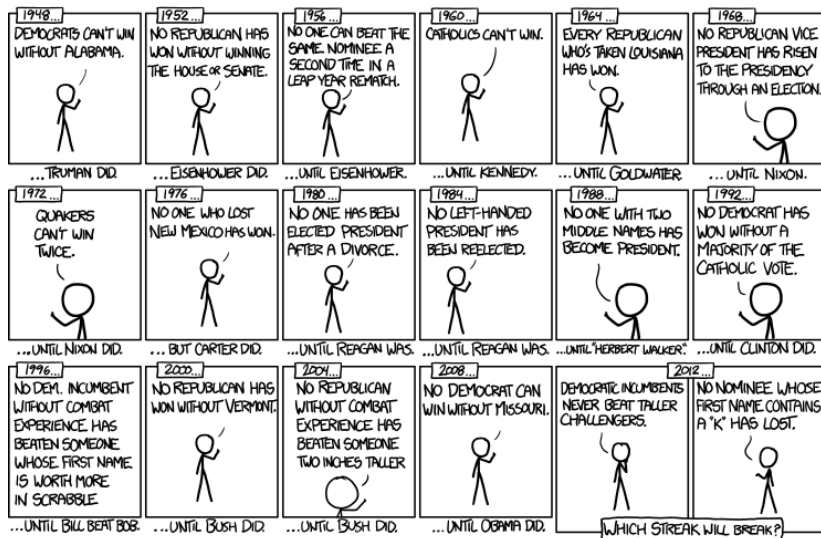
# Complexity control: how many neighbors *k*?



**Small k** gives a **complex** boundary, **large k** results in **coarse** averaging
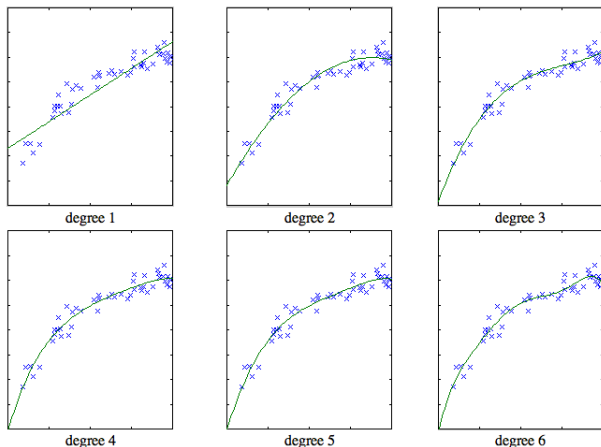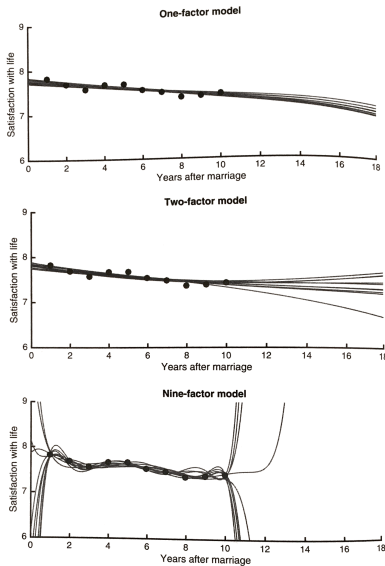
# Overfitting[7]

# Overfitting[7]

# Philosophy

Our models should be complex enough to explain the past, but simple enough to generalize to the future

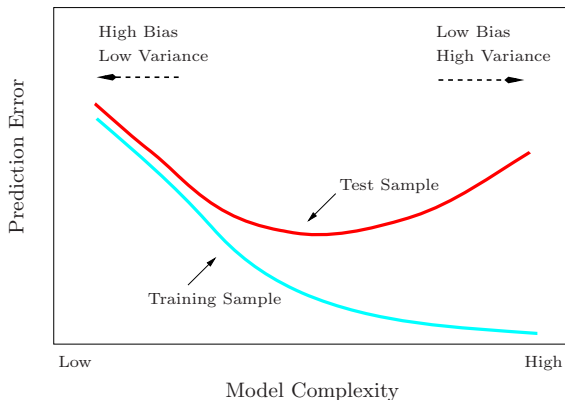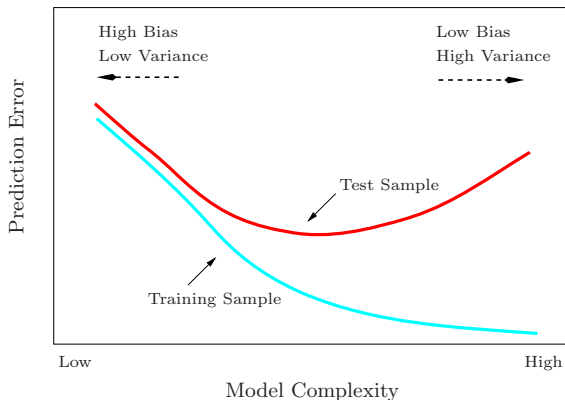# Bias-variance tradeoff

# Bias-variance tradeoff



Simple models may be "wrong" (high bias), but fits don't vary a lot with different samples of training data (low variance)
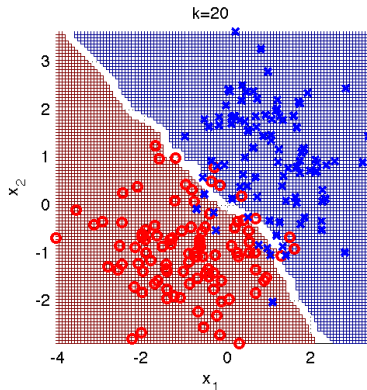
# Bias-variance tradeoff



Flexible models can capture more complex relationships (low bias),
but are also sensitive to noise in the training data (high variance)

# Cross-validation



- Randomly split our data into three sets
- Fit models on the training set
- Use the validation set to find the best model
- Quote final performance of this model on the test set

# Complexity control: how many neighbors $k$?



Evaluate performance on a **held-out test set** to assess
**generalization error**

# Cross-validation for digit recognition



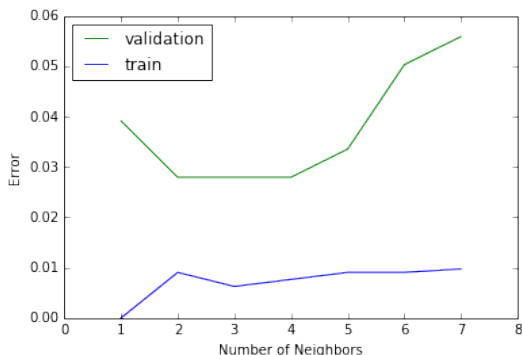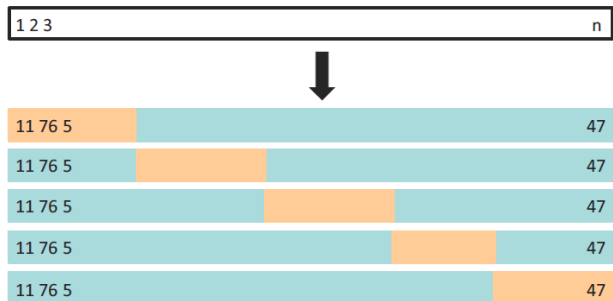- Randomly split our data into three sets
- For each number of neighbors $k$:
  - Fit a model to the training set
  - Evaluate on the validation set
- Select the model with the lowest validation error
- Quote final performance of this model on the test set

# K-fold cross-validation

Estimates of generalization error from one train / validation split can be noisy, so shuffle data and average over $K$ distinct validation partitions instead
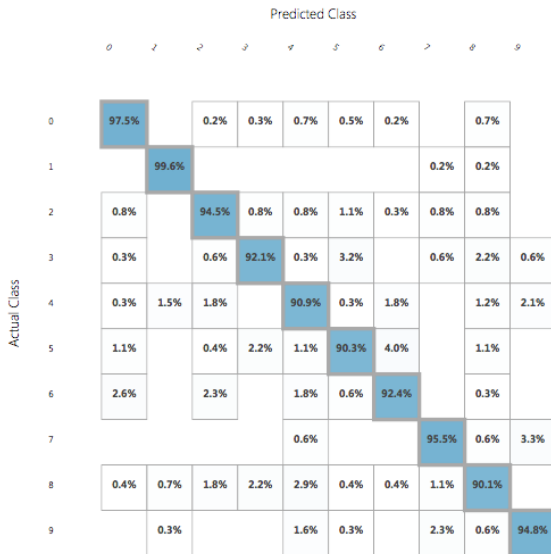
# K-fold cross-validation for digit recognition



- Randomly divide the data into $F$ folds
- For each number of neighbors $k$
    - For each fold:
        - Fit a model on everything but one fold
        - Evaluate the model on the held-out fold
    - Compute the average validation error across folds
- Select the model with the lowest average validation error
- Quote final performance of this model on the test set

# Performance

# k-Nearest Neighbors

Simple approach:
Predict the future by memorizing the past

# k-Nearest Neighbors

Still many choices:

Number of neighbors, feature transformations, distance measures, ...
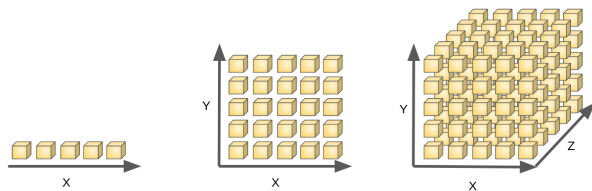
# k-Nearest Neighbors

Problems (practice):

Becomes costly as we see more training examples

# k-Nearest Neighbors

Problems (theory):
Performs poorly when large number of features relative to examples

# Curse of dimensionality[8]



kNN requires a prohibitive number of training examples in high dimensions (e.g., for text data)

# k-Nearest Neighbors

Even with simple methods, there are many decisions to be made in developing a successful machine learning solution

# Bigger models $\neq$ Better models