

Introduction to ML

May 22, 2023

What is Machine Learning?

- Making predictions or decisions from data
- “Programming computers to optimize a performance criterion using example data or past experience” (Ethem Alpaydin, Machine Learning, 2010)
- “A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P, if its performance at tasks in T, as measured by P, improves with experience E.” (Tom Mitchell, Machine Learning, 1997)
- “Learning general models from a data of particular examples”
- “Build a model that is *a good and useful approximation* to the data.”

Today

Traditional Programming



Machine Learning



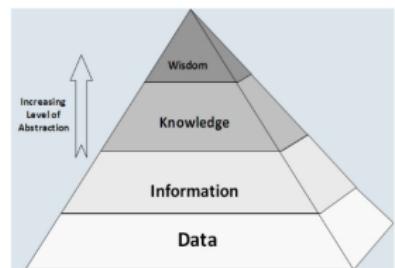
Related Terms

Machine Learning, Data Mining, Knowledge Discovery,
Artificial Intelligence, Statistical Learning, Pattern Recognition,
Computational Learning



When is Machine Learning Used?

- Human expertise does not exist
 - E.g. navigating on Mars
- Humans are unable to explain their expertise
 - E.g. speech recognition
- Solution changes in time
 - E.g. routing on a computer network
- Solution needs to be adapted to particular cases
 - E.g. user biometrics
- Data is cheap and abundant; knowledge is expensive and scarce



Applications of Machine Learning

From: cheapsales@buystufffromme.com
To: ang@cs.stanford.edu
Subject: Buy now!

Deal of the week! Buy now!
Rolex w4tchs - \$100
Medicine (any kind) - \$50
Also low cost M0rgages available.

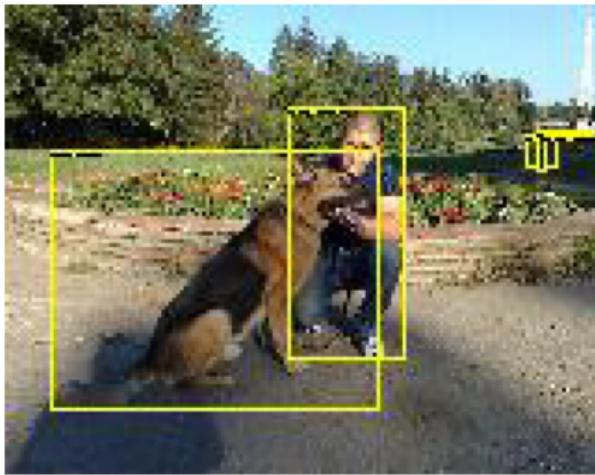
Spam

From: Alfred Ng
To: ang@cs.stanford.edu
Subject: Christmas dates?

Hey Andrew,
Was talking to Mom about plans for Xmas. When do you get off work. Meet Dec 22?
Alf

Non-spam

Applications of Machine Learning



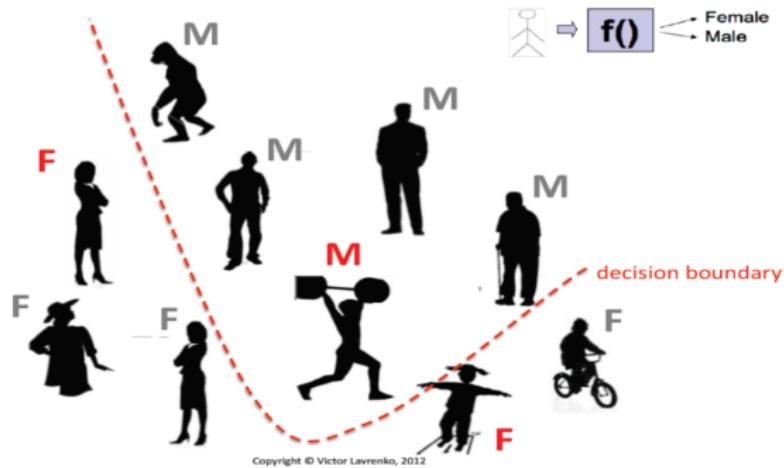
More ML Applications

- Science (Astronomy, neuroscience, medical imaging, bio-informatics)
- Environment (energy, climate, weather, resources)
- Retail (Intelligent stock control, demographic store placement)
- Manufacturing (Intelligent control, automated monitoring, detection methods)
- Security (Intelligent smoke alarms, fraud detection)
- Marketing (promotions, ...)
- Management (Scheduling, timetabling)
- Finance (credit scoring, risk analysis...)
- Web data (information retrieval, information extraction, ...)

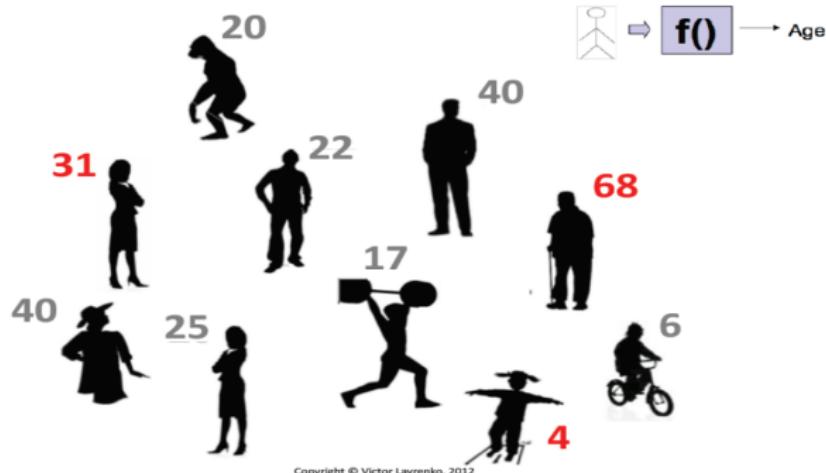
Overview of ML

- Supervised learning
 - Predict an output y when given an input x
 - For categorical y : classification.
 - For real-valued y : regression.
- Unsupervised learning
 - Create an internal representation of the input, e.g. clustering, dimensionality
 - This is important in machine learning as getting labels is often difficult and expensive
- Other settings of ML
 - Reinforcement learning (learning from “rewards”)
 - Semi-supervised learning (combines supervised + unsupervised)
 - Active learning, Transfer learning, Structured prediction

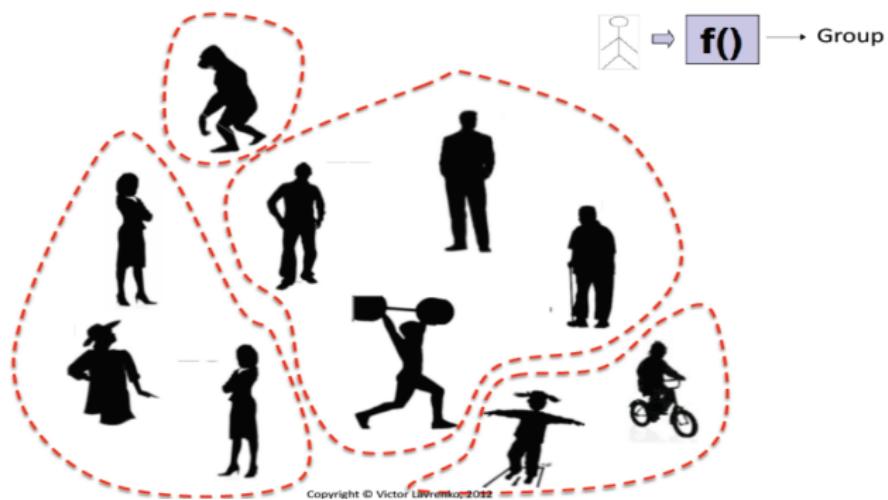
Classification (Supervised Learning)



Regression (Supervised Learning)

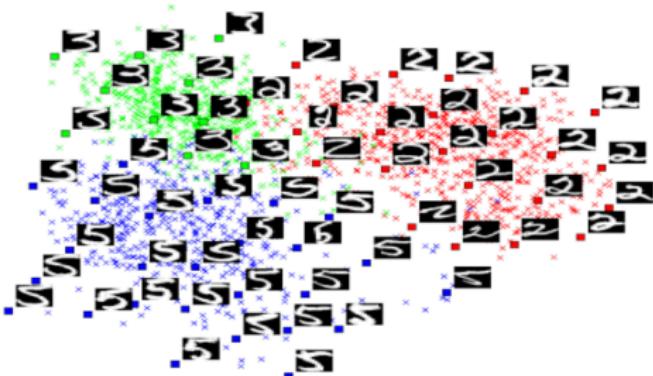


Clustering (Unsupervised Learning)



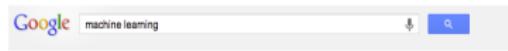
Dimensionality Reduction (Unsupervised Learning)

- Large sample size is required for high-dimensional data
- Query accuracy and efficiency degrade rapidly as the dimension increases
- Strategies
 - Feature reduction
 - Feature selection
 - Manifold learning
 - Kernel learning



Other Settings: Ranking (Supervised Learning)

Given a query and
a set of web pages,
rank them according
to relevance



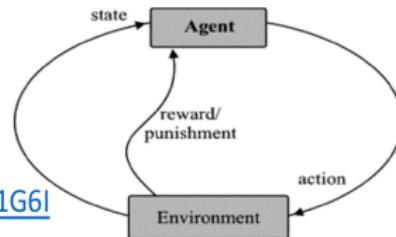
- About 130,000,000 results (0.26 seconds)
- Machine learning - Wikipedia, the free encyclopedia**
en.wikipedia.org/wiki/Machine_learning •
Machine learning, a branch of artificial intelligence, concerns the construction and study of systems that can learn from data. For example, a machine learning ...
 - Artificial Intelligence - Supervised learning - List of machine learning ... - Webs**
Franck Demenoult +1d this
 - CS 229, Machine Learning**
cs229.stanford.edu/ •
Check out this year's awesome projects at Fall 2012 Projects. Come check out the cool new projects during the CS229 Poster Session this Thursday December ...
You've visited this page 2 times. Last visit: 8/14/13
 - Machine Learning | Coursera**
<https://www.coursera.org/courses/ml> •
Machine learning is the science of getting computers to act without being explicitly programmed. In the past decade, machine learning has given us self-driving ...
 - Machine Learning Department - Carnegie Mellon University**
www.ml.cmu.edu/ •
Large group with projects in robot learning, data mining for manufacturing and in multimedia databases, causal inference, and disclosure limitation.
 - Machine Learning - MIT OpenCourseWare**
ocw.mit.edu/courses/Electrical-Engineering-and-Computer-Science/6.S071/ •
6.S071 is an introductory course on machine learning which gives an overview of many concepts, techniques, and algorithms in machine learning, beginning with ...

• Other applications

- User preference, e.g. Netflix “My List” -- movie queue ranking
- Flight search (search in general)
- ...

Other Settings: Reinforcement Learning

- Learning a policy: A **sequence** of outputs
- No supervised output but delayed reward
 - E.g. Game playing
 - E.g. Robot in a maze
- Multiple agents, partial observability, ...
- Example (Simple Demo):
 - <https://www.youtube.com/watch?v=DCjbk4m1G6I>



ML Problems

| | <i>Supervised Learning</i> | <i>Unsupervised Learning</i> |
|-------------------|----------------------------------|------------------------------|
| <i>Discrete</i> | classification or categorization | clustering |
| <i>Continuous</i> | regression | dimensionality reduction |

Mathematical Basis

- Functions, Logarithms and Exponentials
- Vectors, Dot Products, Orthogonality
- Matrices, Matrix Operations, Linear Transformations, Eigendecomposition
- Calculus, Differentiation, Integration
- Probability and Statistics
- Functional Analysis, Hilbert Spaces

ML Problems: Recall

| | <i>Supervised Learning</i> | <i>Unsupervised Learning</i> |
|-------------------|----------------------------------|------------------------------|
| <i>Discrete</i> | classification or categorization | clustering |
| <i>Continuous</i> | regression | dimensionality reduction |

Classification Methods

- k-Nearest Neighbors
- Decision Trees
- Naïve Bayes
- Support Vector Machines
- Logistic Regression
- Neural Networks
- Ensemble Methods (Boosting, Random Forests)

How to evaluate?

Training vs Generalization Error

- Training Error
 - Not very useful
 - Relatively easy to obtain low error
- Generalization Error
 - How well we do on future data

How to compute
generalization error?

$$E_{train} = \frac{1}{n} \sum_{i=1}^n \underbrace{\text{error}}_{\substack{\text{same? different by how much?} \\ \text{training examples}}}(\underbrace{f_D(\mathbf{x}_i)}_{\substack{\text{value we predicted}}}, \underbrace{y_i}_{\substack{\text{true value}}})$$

$$E_{gen} = \int \underbrace{\text{error}}_{\substack{\text{over all possible } x,y \\ \text{error as before}}}(f_D(\mathbf{x}), y) p(y, \mathbf{x}) d\mathbf{x}$$

how often we expect to see such x and y

Estimating Generalization Error

- Testing Error

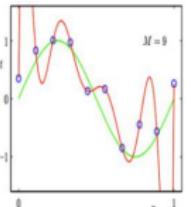
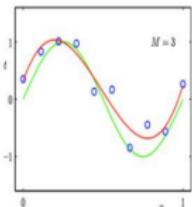
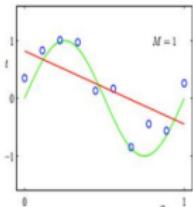
- Set aside part of training data (testing set)
- Learn a predictor without using any of this test data
- Predict values for testing set, compute error
- This is an estimate of generalization error

$$E_{test} = \frac{1}{n} \sum_{i=1}^n \text{error}(f_D(\mathbf{x}_i), y_i)$$

over testing set

Underfitting and Overfitting

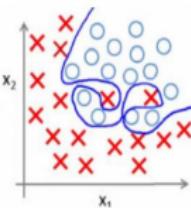
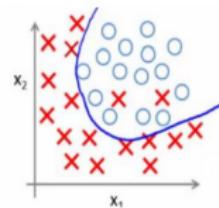
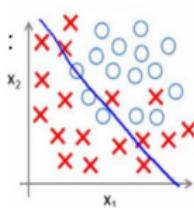
Regression



predictor too inflexible:
cannot capture pattern

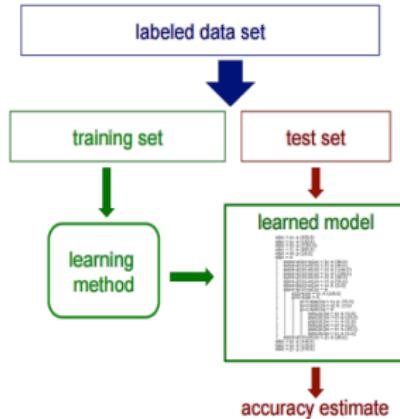
predictor too flexible:
fits noise in the data

Classification

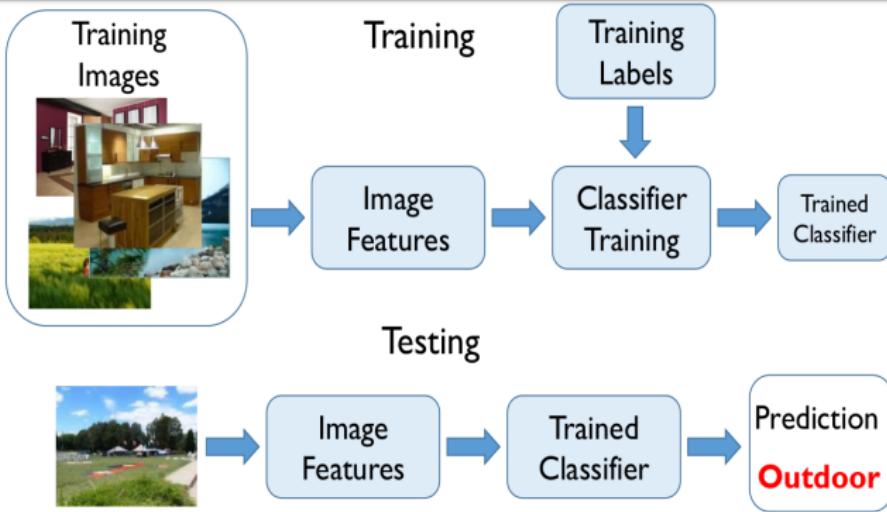


Estimating Generalization Error

- Getting an unbiased estimate of the accuracy of a learned model



Example: Image Classification



Training, Validation, Test Sets

Training set

- NB: Count frequencies, DT: Pick attributes to split on

Validation set

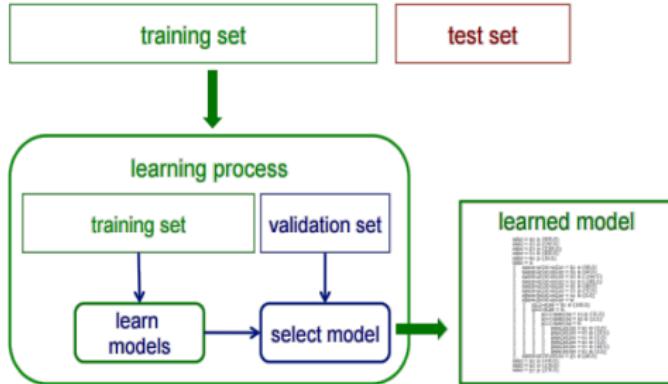
- Pick best-performing algorithm (NB vs DT vs..)
- Fine-tune parameters (Tree depth, k in kNN, c in SVM)

Testing set

- Run multiple trials and average

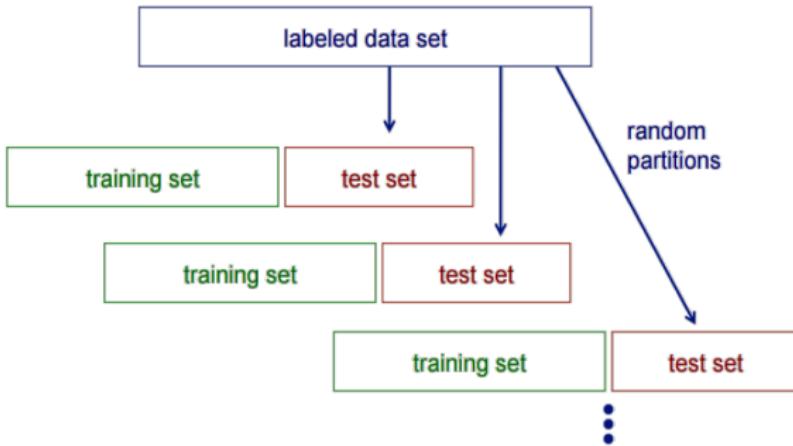
Use of Validation Sets

- If we want unbiased estimates of accuracy during the learning process:



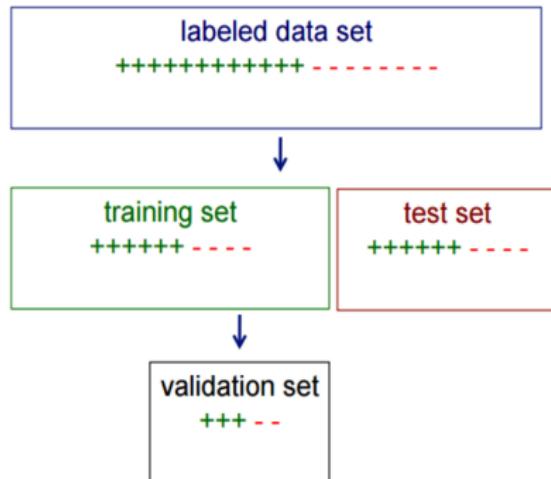
Random Resampling

- We can artificially increase training set size using **random resampling**:



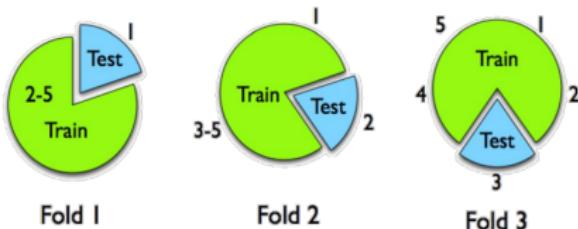
Stratified Sampling

- When randomly selecting training or validation sets, we may want to ensure that class proportions are maintained in each selected set
- This can be done via **stratified sampling**: first stratify instances by class, then randomly select instances from each class proportionally.



Model Selection

- Resubstitution
- K-fold cross-validation



- Leave-one-out
 - N-fold cross-validation

Cross-Validation: Example

- Suppose we have 100 instances, and we want to estimate accuracy with cross validation

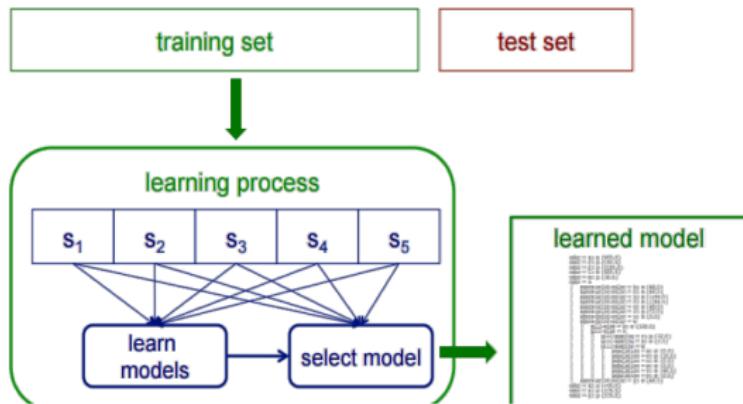
| iteration | train on | test on | correct |
|-----------|---|----------------|---------|
| 1 | s ₂ s ₃ s ₄ s ₅ | s ₁ | 11 / 20 |
| 2 | s ₁ s ₃ s ₄ s ₅ | s ₂ | 17 / 20 |
| 3 | s ₁ s ₂ s ₄ s ₅ | s ₃ | 16 / 20 |
| 4 | s ₁ s ₂ s ₃ s ₅ | s ₄ | 13 / 20 |
| 5 | s ₁ s ₂ s ₃ s ₄ | s ₅ | 16 / 20 |

$$\text{Classification Accuracy} = 73/100 = 73\%$$

Note: Whenever we use multiple training sets, as in CV and random resampling, we are evaluating a learning method as opposed to an individual learned model

Cross-Validation: Example

- Instead of a single validation set, we can use cross-validation within a training set to select a model (e.g. to choose the best k in k-NN)



Evaluation Measures

- Classification
 - How often we classify something right/wrong
- Regression
 - How close are we to what we're trying to predict
- Ranking/Search
 - How correct are the top-k results?
- Clustering
 - How well we describe our data (Not straightforward)

Is accuracy adequate?

- Accuracy may not be useful in cases where
 - There is a large class skew
 - Is 98% accuracy good if 97% of the instances are negative?
 - There are differential misclassification costs – say, getting a positive wrong costs more than getting a negative wrong
 - Consider a medical domain in which a false positive results in an extraneous test but a false negative results in a failure to treat a disease
 - We are most interested in a subset of high-confidence predictions

Classification Error: Beyond Accuracy

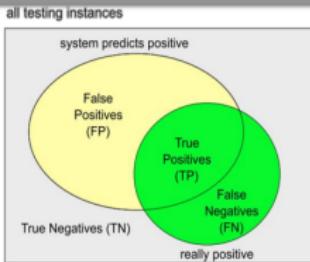
In 2-class problems:

| | | actual class | |
|-----------------|----------|----------------------|----------------------|
| | | positive | negative |
| predicted class | positive | true positives (TP) | false positives (FP) |
| | negative | false negatives (FN) | true negatives (TN) |

$$\text{accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{FN} + \text{TN}}$$

Classification Performance Measures

| | | Predict positive? | |
|------------------|-----|-------------------|----|
| | | Yes | No |
| Really positive? | Yes | TP | FN |
| | No | FP | TN |



- Classification Error: $\frac{errors}{total} = \frac{FP+FN}{TP+TN+FP+FN}$
- Accuracy = 1-Error: $\frac{correct}{total} = \frac{TP+TN}{TP+TN+FP+FN}$
- False Alarm = False Positive rate = $FP / (FP+TN)$
- Miss = False Negative rate = $FN / (TP+FN)$
- Recall = True Positive rate = $TP / (TP+FN)$
- Precision = $TP / (TP+FP)$

meaningless
if classes
imbalanced

always report
in pairs, e.g.:
Miss / FA or
Recall / Prec.

- “Sensitivity” = Probability of a positive test given a patient has the disease
- “Specificity” = Probability of a negative test given a patient is well

Classification Error: Beyond Accuracy

For multi-class problems?

Confusion Matrix

