

CS 188: Artificial Intelligence

Spring 2011

Lecture 21: Perceptrons

4/13/2010

Pieter Abbeel – UC Berkeley
Many slides adapted from Dan Klein.

Announcements

- Project 4: due Friday.
- Final Contest: up and running!
- Project 5 out!
- Saturday, 10am-noon, 3rd floor
Sutardja Dai Hall

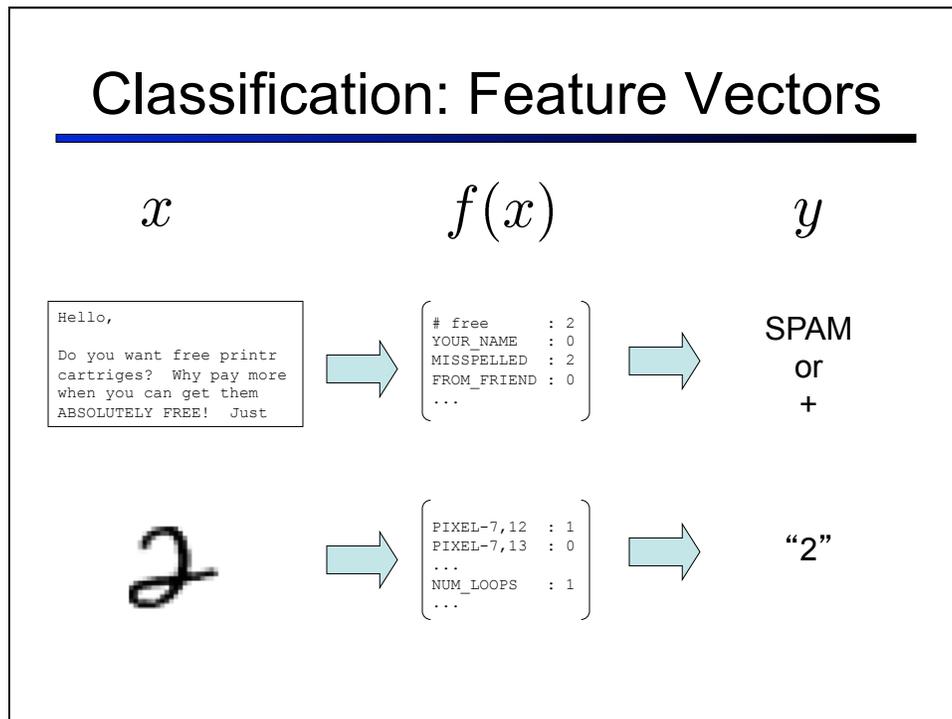


Survey

Outline

- Generative vs. Discriminative
- Perceptron

Classification: Feature Vectors



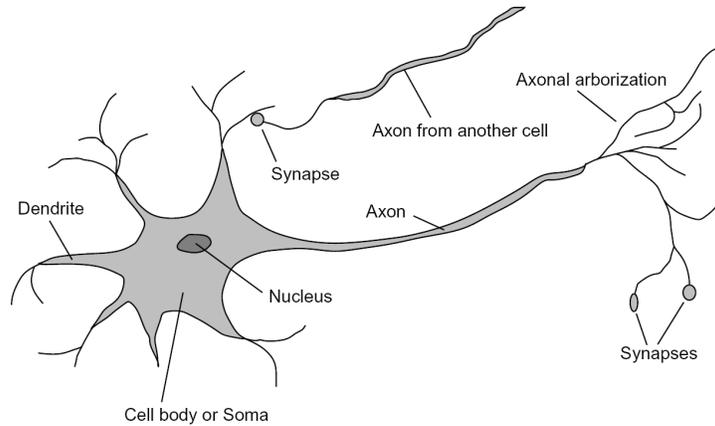
Generative vs. Discriminative

- **Generative classifiers:**
 - E.g. naïve Bayes
 - A causal model with evidence variables
 - Query model for causes given evidence
- **Discriminative classifiers:**
 - No causal model, no Bayes rule, often no probabilities at all!
 - Try to predict the label Y directly from X
 - Robust, accurate with varied features
 - Loosely: **mistake driven rather than model driven**

6

Some (Simplified) Biology

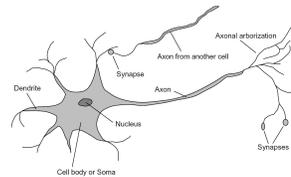
- Very loose inspiration: human neurons



7

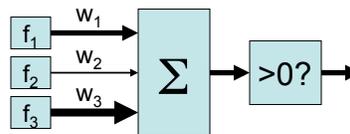
Linear Classifiers

- Inputs are **feature values**
- Each feature has a **weight**
- Sum is the **activation**



$$\text{activation}_w(x) = \sum_i w_i \cdot f_i(x) = w \cdot f(x)$$

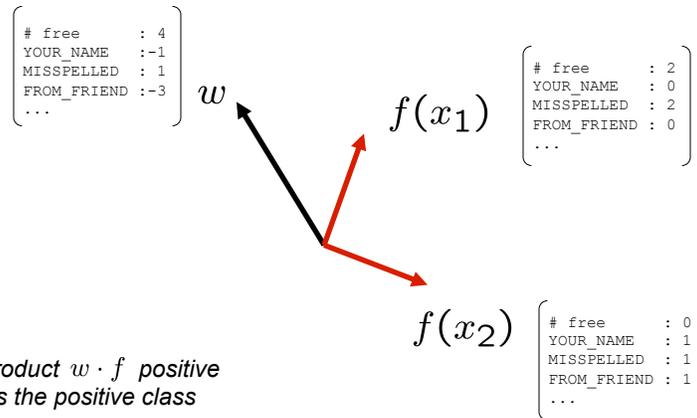
- If the activation is:
 - Positive, output +1
 - Negative, output -1



8

Classification: Weights

- Binary case: compare features to a weight vector
- Learning: figure out the weight vector from examples

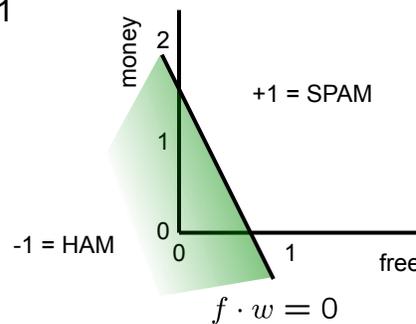


Binary Decision Rule

- In the space of feature vectors
 - Examples are points
 - Any weight vector is a hyperplane
 - One side corresponds to $Y=+1$
 - Other corresponds to $Y=-1$

w

| | |
|-------|------|
| BIAS | : -3 |
| free | : 4 |
| money | : 2 |
| ... | |



Outline

- Naïve Bayes recap
- Smoothing
- Generative vs. Discriminative
- *Perceptron*

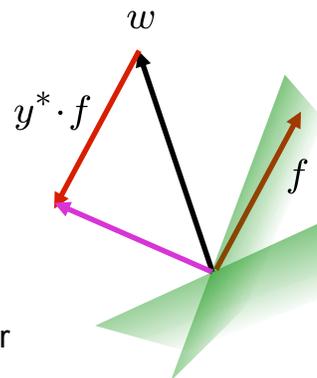
Binary Perceptron Update

- Start with zero weights
- For each training instance:
 - Classify with current weights

$$y = \begin{cases} +1 & \text{if } w \cdot f(x) \geq 0 \\ -1 & \text{if } w \cdot f(x) < 0 \end{cases}$$

- If correct (i.e., $y=y^*$), no change!
- If wrong: adjust the weight vector by adding or subtracting the feature vector. Subtract if y^* is -1.

$$w = w + y^* \cdot f$$



[demo] 14

Multiclass Decision Rule

- If we have multiple classes:
 - A weight vector for each class:

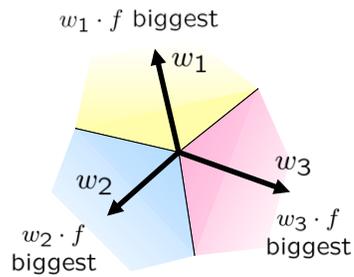
$$w_y$$

- Score (activation) of a class y:

$$w_y \cdot f(x)$$

- Prediction highest score wins

$$y = \arg \max_y w_y \cdot f(x)$$



Binary = multiclass where the negative class has weight zero

Example

“win the vote” →

| | |
|------|-----|
| BIAS | : 1 |
| win | : 1 |
| game | : 0 |
| vote | : 1 |
| the | : 1 |
| ... | |

w_{SPORTS}

| | |
|------|------|
| BIAS | : -2 |
| win | : 4 |
| game | : 4 |
| vote | : 0 |
| the | : 0 |
| ... | |

$w_{POLITICS}$

| | |
|------|-----|
| BIAS | : 1 |
| win | : 2 |
| game | : 0 |
| vote | : 4 |
| the | : 0 |
| ... | |

w_{TECH}

| | |
|------|-----|
| BIAS | : 2 |
| win | : 0 |
| game | : 2 |
| vote | : 0 |
| the | : 0 |
| ... | |

Learning Multiclass Perceptron

- Start with zero weights
- Pick up training instances one by one
- Classify with current weights

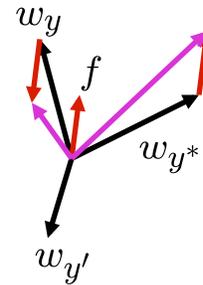
$$y = \arg \max_y w_y \cdot f(x)$$

$$= \arg \max_y \sum_i w_{y,i} \cdot f_i(x)$$

- If correct, no change!
- If wrong: lower score of wrong answer, raise score of right answer

$$w_y = w_y - f(x)$$

$$w_{y^*} = w_{y^*} + f(x)$$



17

Example

“win the vote”

“win the election”

“win the game”

w_{SPORTS}

| | |
|------|---|
| BIAS | : |
| win | : |
| game | : |
| vote | : |
| the | : |
| ... | : |

$w_{POLITICS}$

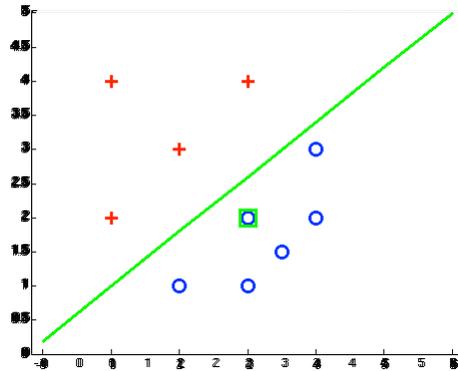
| | |
|------|---|
| BIAS | : |
| win | : |
| game | : |
| vote | : |
| the | : |
| ... | : |

w_{TECH}

| | |
|------|---|
| BIAS | : |
| win | : |
| game | : |
| vote | : |
| the | : |
| ... | : |

Examples: Perceptron

- Separable Case



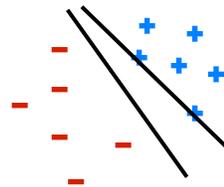
19

Properties of Perceptrons

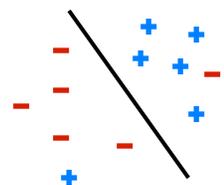
- Separability: some parameters get the training set perfectly correct
- Convergence: if the training is separable, perceptron will eventually converge (binary case)
- Mistake Bound: the maximum number of mistakes (binary case) related to the *margin* or degree of separability

$$\text{mistakes} < \frac{k}{\delta^2}$$

Separable



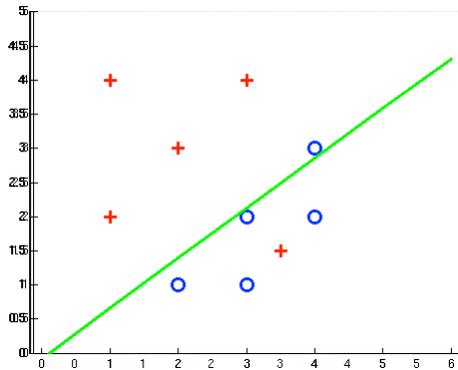
Non-Separable



21

Examples: Perceptron

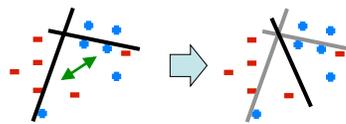
Non-Separable Case



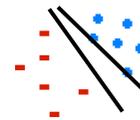
22

Problems with the Perceptron

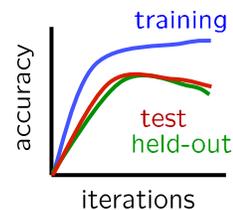
- Noise: if the data isn't separable, weights might thrash
 - Averaging weight vectors over time can help (averaged perceptron)



- Mediocre generalization: finds a "barely" separating solution

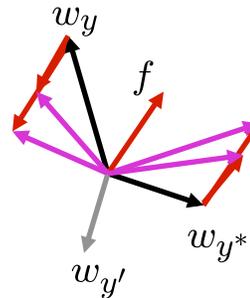


- Overtraining: test / held-out accuracy usually rises, then falls
 - Overtraining is a kind of overfitting



Fixing the Perceptron

- Idea: adjust the weight update to mitigate these effects
- MIRA*: choose an update size that fixes the current mistake...
- ... but, minimizes the change to w



$$\min_w \frac{1}{2} \sum_y \|w_y - w'_y\|^2$$

$$w_{y^*} \cdot f(x) \geq w_y \cdot f(x) + 1$$

- The +1 helps to generalize

* Margin Infused Relaxed Algorithm

Guessed y instead of y^* on example x with features $f(x)$

$$w_y = w'_y - \tau f(x)$$

$$w_{y^*} = w'_{y^*} + \tau f(x)$$

Minimum Correcting Update

$$\min_w \frac{1}{2} \sum_y \|w_y - w'_y\|^2$$

$$w_{y^*} \cdot f \geq w_y \cdot f + 1$$



$$\min_{\tau} \|\tau f\|^2$$

$$w_{y^*} \cdot f \geq w_y \cdot f + 1$$



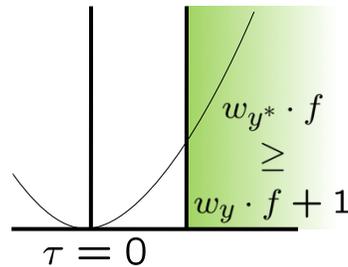
$$\min_{\tau} \tau^2$$

$$(w'_{y^*} + \tau f) \cdot f \geq (w'_y - \tau f) \cdot f + 1$$

$$\tau \geq \frac{(w'_y - w'_{y^*}) \cdot f + 1}{2f \cdot f}$$

$$w_y = w'_y - \tau f(x)$$

$$w_{y^*} = w'_{y^*} + \tau f(x)$$



min not $\tau=0$, or would not have made an error, so min will be where equality holds

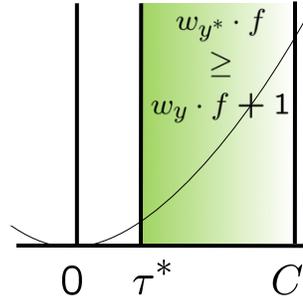
Maximum Step Size

- In practice, it's also bad to make updates that are too large

- Example may be labeled incorrectly
- You may not have enough features
- Solution: cap the maximum possible value of τ with some constant C

$$\tau^* = \min \left(\frac{(w'_y - w'_{y^*}) \cdot f + 1}{2f \cdot f}, C \right)$$

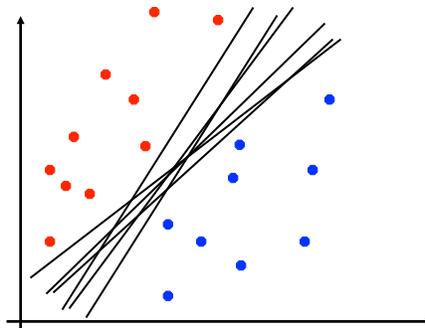
- Corresponds to an optimization that assumes non-separable data
- Usually converges faster than perceptron
- Usually better, especially on noisy data



27

Linear Separators

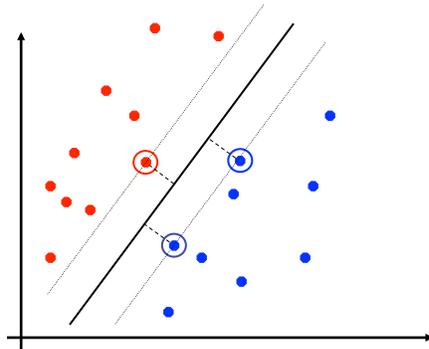
- Which of these linear separators is optimal?



28

Support Vector Machines

- **Maximizing the margin:** good according to intuition, theory, practice
- Only **support vectors** matter; other training examples are ignorable
- Support vector machines (SVMs) find the separator with max margin
- Basically, SVMs are MIRA where you optimize over all examples at once



MIRA

$$\min_w \frac{1}{2} \|w - w'\|^2$$
$$w_{y^*} \cdot f(x_i) \geq w_y \cdot f(x_i) + 1$$

SVM

$$\min_w \frac{1}{2} \|w\|^2$$
$$\forall i, y \quad w_{y^*} \cdot f(x_i) \geq w_y \cdot f(x_i) + 1$$

Classification: Comparison

- **Naïve Bayes**
 - Builds a model training data
 - Gives prediction probabilities
 - Strong assumptions about feature independence
 - One pass through data (counting)
- **Perceptrons / MIRA:**
 - Makes less assumptions about data
 - Mistake-driven learning
 - Multiple passes through data (prediction)
 - Often more accurate

30

Extension: Web Search

- Information retrieval:
 - Given information needs, produce information
 - Includes, e.g. web search, question answering, and classic IR
- Web search: not exactly classification, but rather ranking

$x = \text{“Apple Computers”}$



Feature-Based Ranking

$x = \text{“Apple Computers”}$

$$f(x, \text{Screenshot of 'Apple' search result}) = [0.3 \ 5 \ 0 \ 0 \ \dots]$$

$$f(x, \text{Screenshot of 'Apple Inc.' search result}) = [0.8 \ 4 \ 2 \ 1 \ \dots]$$

Perceptron for Ranking

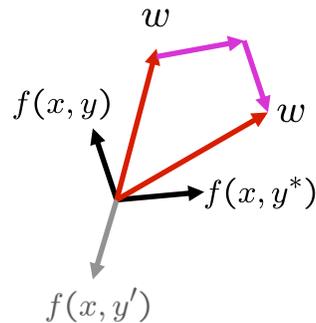
- Inputs x
- Candidates y
- Many feature vectors: $f(x, y)$
- One weight vector: w

- Prediction:

$$y = \arg \max_y w \cdot f(x, y)$$

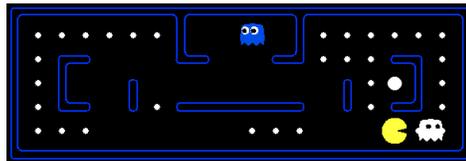
- Update (if wrong):

$$w = w + f(x, y^*) - f(x, y)$$



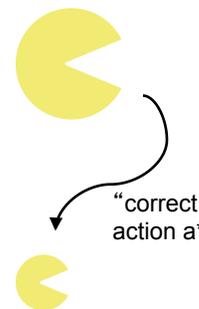
Pacman Apprenticeship!

- Examples are states s



- Candidates are pairs (s, a)
- “Correct” actions: those taken by expert
- Features defined over (s, a) pairs: $f(s, a)$
- Score of a q -state (s, a) given by:

$$w \cdot f(s, a)$$



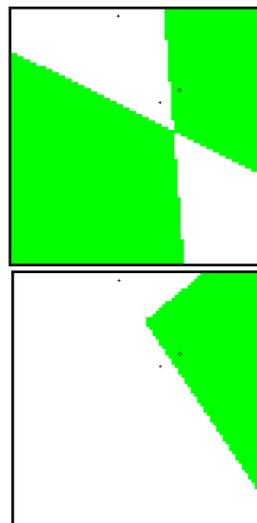
$$\forall a \neq a^*, w \cdot f(a^*) > w \cdot f(a)$$

- How is this VERY different from reinforcement learning?



Case-Based Reasoning

- **Similarity for classification**
 - Case-based reasoning
 - Predict an instance's label using similar instances
- **Nearest-neighbor classification**
 - 1-NN: copy the label of the most similar data point
 - K-NN: let the k nearest neighbors vote (have to devise a weighting scheme)
 - Key issue: how to define similarity
 - Trade-off:
 - Small k gives relevant neighbors
 - Large k gives smoother functions
 - Sound familiar?
- [Demo]

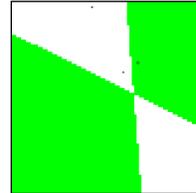


36

<http://www.cs.cmu.edu/~zhuxj/courseproject/knndemo/KNN.html>

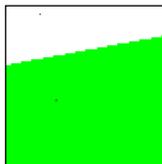
Parametric / Non-parametric

- **Parametric models:**
 - Fixed set of parameters
 - More data means better settings
- **Non-parametric models:**
 - Complexity of the classifier increases with data
 - Better in the limit, often worse in the non-limit
- (K)NN is **non-parametric**

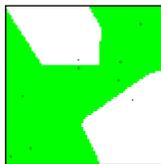


Truth

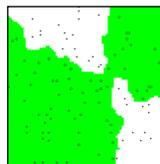
2 Examples



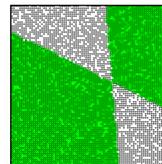
10 Examples



100 Examples



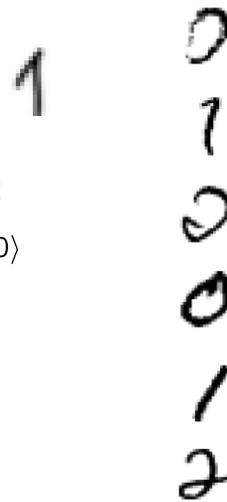
10000 Examples



37

Nearest-Neighbor Classification

- **Nearest neighbor for digits:**
 - Take new image
 - Compare to all training images
 - Assign based on closest example



- **Encoding: image is vector of intensities:**

$$1 = \langle 0.0 \ 0.0 \ 0.3 \ 0.8 \ 0.7 \ 0.1 \ \dots \ 0.0 \rangle$$

- **What's the similarity function?**
 - Dot product of two images vectors?

$$\text{sim}(x, x') = x \cdot x' = \sum_i x_i x'_i$$

- Usually normalize vectors so $\|x\| = 1$
- min = 0 (when?), max = 1 (when?)

38

Basic Similarity

- Many similarities based on **feature dot products**:

$$\text{sim}(x, x') = f(x) \cdot f(x') = \sum_i f_i(x) f_i(x')$$

- If features are just the pixels:

$$\text{sim}(x, x') = x \cdot x' = \sum_i x_i x'_i$$

- Note: not all similarities are of this form

39

Invariant Metrics

- Better distances use knowledge about vision
- Invariant metrics:
 - Similarities are invariant under certain transformations
 - Rotation, scaling, translation, stroke-thickness...
 - E.g: 
 - 16 x 16 = 256 pixels; a point in 256-dim space
 - Small similarity in \mathbb{R}^{256} (why?)
 - How to incorporate invariance into similarities?

40

This and next few slides adapted from Xiao Hu, UIUC

Template Deformation

- **Deformable templates:**
 - An “ideal” version of each category
 - Best-fit to image using min variance
 - Cost for high distortion of template
 - Cost for image points being far from distorted template
- **Used in many commercial digit recognizers**



43
Examples from [Hastie 94]

A Tale of Two Approaches...

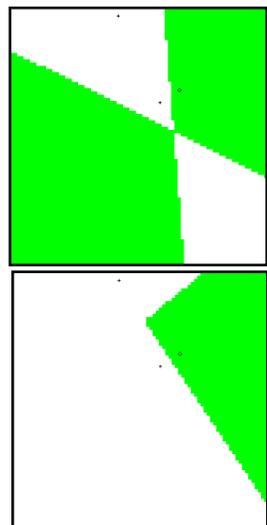
- **Nearest neighbor-like approaches**
 - Can use fancy similarity functions
 - Don't actually get to do explicit learning
- **Perceptron-like approaches**
 - Explicit training to reduce empirical error
 - Can't use fancy similarity, only linear
 - Or can they? Let's find out!

44



Case-Based Reasoning

- **Similarity for classification**
 - Case-based reasoning
 - Predict an instance's label using similar instances
- **Nearest-neighbor classification**
 - 1-NN: copy the label of the most similar data point
 - K-NN: let the k nearest neighbors vote (have to devise a weighting scheme)
 - Key issue: how to define similarity
 - Trade-off:
 - Small k gives relevant neighbors
 - Large k gives smoother functions
 - Sound familiar?
- [DEMO]

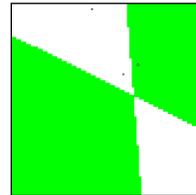


<http://www.cs.cmu.edu/~zhuxj/courseproject/knndemo/KNN.html>

Recap: Nearest-Neighbor

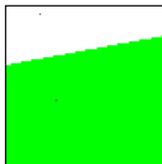
- **Nearest neighbor:**

- Classify test example based on closest training example
- Requires a similarity function (**kernel**)
- **Eager learning:** extract classifier from data
- **Lazy learning:** keep data around and predict from it at test time

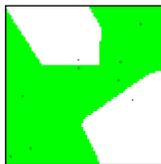


Truth

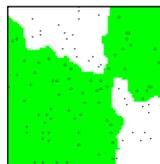
2 Examples



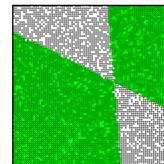
10 Examples



100 Examples



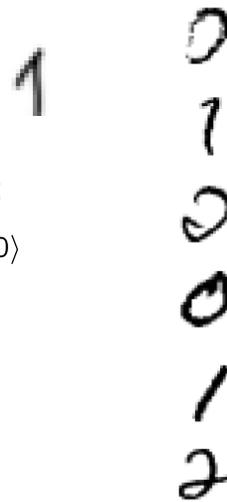
10000 Examples



Nearest-Neighbor Classification

- **Nearest neighbor for digits:**

- Take new image
- Compare to all training images
- Assign based on closest example



- **Encoding: image is vector of intensities:**

$$1 = \langle 0.0 \ 0.0 \ 0.3 \ 0.8 \ 0.7 \ 0.1 \ \dots \ 0.0 \rangle$$

- **What's the similarity function?**

- Dot product of two images vectors?

$$\text{sim}(x, y) = x \cdot y = \sum_i x_i y_i$$

- Usually normalize vectors so $\|x\| = 1$
- min = 0 (when?), max = 1 (when?)