• Features can be much more complex

- Features can be much more complex
- Drawn from bigger discrete set

- Features can be much more complex
- Drawn from bigger discrete set
 - If set is unordered (4 different makes of cars, for example), use binary attributes to encode the values (1000, 0100, 0010, 0001)

- Features can be much more complex
- Drawn from bigger discrete set
 - If set is unordered (4 different makes of cars, for example), use binary attributes to encode the values (1000, 0100, 0010, 0001)
 - If set is ordered, treat as real-valued

- Features can be much more complex
- Drawn from bigger discrete set
 - If set is unordered (4 different makes of cars, for example), use binary attributes to encode the values (1000, 0100, 0010, 0001)
 - If set is ordered, treat as real-valued
- Real-valued: bias that inputs whose features have "nearby" values ought to have "nearby" outputs



Love thy Nearest Neighbor

- Remember all your data
- When someone asks a question,
 - -find the nearest old data point
 - -return the answer associated with it



What do we mean by "Nearest"?

- Need a distance function on inputs
- Typically use Euclidean distance (length of a straight line between the points)

$$D(x^i, x^k) = \sqrt{\sum_j (x^i_j - x^k_j)^2}$$

What do we mean by "Nearest"?

- Need a distance function on inputs
- Typically use Euclidean distance (length of a straight line between the points)

$$D(x^{i}, x^{k}) = \sqrt{\sum_{j} (x_{j}^{i} - x_{j}^{k})^{2}}$$

 Distance between character strings might be number of edits required to turn one into the other

- What if we're trying to predict a car's gas mileage?
 - f_1 = weight in pounds
 - f_2 = number of cylinders

- What if we're trying to predict a car's gas mileage?
 - f_1 = weight in pounds
 - f_2 = number of cylinders
- Any effect of f_2 will be completely lost because of the relative scales

- What if we're trying to predict a car's gas mileage?
 - f_1 = weight in pounds
 - f_2 = number of cylinders
- Any effect of f_2 will be completely lost because of the relative scales
- So, re-scale the inputs

- What if we're trying to predict a car's gas mileage?
 - f_1 = weight in pounds
 - f_2 = number of cylinders
- Any effect of f_2 will be completely lost because of the relative scales
- So, re-scale the inputs to have mean 0 and variance 1:

$$x' = \frac{x - \overline{x}}{\sigma_x} \text{ average}$$

- What if we're trying to predict a car's gas mileage?
 - f_1 = weight in pounds
 - f_2 = number of cylinders
- Any effect of f_2 will be completely lost because of the relative scales
- So, re-scale the inputs to have mean 0 and variance 1:



Or, build knowledge in by scaling features differently

- What if we're trying to predict a car's gas mileage?
 - f_1 = weight in pounds
 - f_2 = number of cylinders
- Any effect of f_2 will be completely lost because of the relative scales
- So, re-scale the inputs to have mean 0 and variance 1:



- Or, build knowledge in by scaling features differently
- Or use cross-validation to choose scales



6.034 - Spring 03 • 16









6.034 - Spring 03 • 20





Hypothesis



6.034 - Spring 03 • 23

Hypothesis



6.034 - Spring 03 • 24

• Learning is fast

- Learning is fast
- Lookup takes about m*n computations

- Learning is fast
- Lookup takes about m*n computations
 - storing data in a clever data structure (KD-tree) reduces this, on average, to log(m)*n

- Learning is fast
- Lookup takes about m*n computations
 - storing data in a clever data structure (KD-tree) reduces this, on average, to log(m)*n
- Memory can fill up with all that data

- Learning is fast
- Lookup takes about m*n computations
 - storing data in a clever data structure (KD-tree) reduces this, on average, to log(m)*n
- Memory can fill up with all that data
 - delete points that are far away from the boundary

Noise



Noise



Noise



k-Nearest Neighbor



• Find the k nearest points

k-Nearest Neighbor



- Find the k nearest points
- Predict output according to the majority

k-Nearest Neighbor



- Find the k nearest points
- Predict output according to the majority
- Choose k using cross-validation

Curse of Dimensionality

- Nearest neighbor is great in low dimensions (up to about 6)
- As n increases, things get weird:
Curse of Dimensionality

- Nearest neighbor is great in low dimensions (up to about 6)
- As n increases, things get weird:
 - In high dimensions, almost all points are far away from one another
 - They' re almost all near the boundaries

Curse of Dimensionality

- Nearest neighbor is great in low dimensions (up to about 6)
- As n increases, things get weird:
 - In high dimensions, almost all points are far away from one another
 - They' re almost all near the boundaries
- Imagine sprinkling data points uniformly within a 10-dimensional unit cube
 - To capture 10% of the points, you'd need a cube with sides of length .63!

Curse of Dimensionality

- Nearest neighbor is great in low dimensions (up to about 6)
- As n increases, things get weird:
 - In high dimensions, almost all points are far away from one another
 - They' re almost all near the boundaries
- Imagine sprinkling data points uniformly within a 10-dimensional unit cube
 - To capture 10% of the points, you'd need a cube with sides of length .63!
- Cure: feature selection or more global models

Test Domains

6.034 - Spring 03 • 40

Test Domains

- Heart Disease: predict whether a person has significant narrowing of the arteries, based on tests
 - 26 features
 - 297 data points

Test Domains

- Heart Disease: predict whether a person has significant narrowing of the arteries, based on tests
 - 26 features
 - 297 data points

- Auto MPG: predict whether a car gets more than 22 miles per gallon, based on attributes of car
 - 12 features
 - 385 data points

Heart Disease

• Relatively insensitive to k



Heart Disease

- Relatively insensitive to k
- Normalization matters!



Auto MPG

- Relatively insensitive to k
- Normalization doesn't matter much



Auto MPG

- Now normalization matters a lot!
- Watch the scales on your graphs



Remember Decision Trees

Use all the data to build a tree of questions with answers at the leaves



6.034 - Spring 03 • 47

• Tests in nodes can be of the form $x_j > constant$

- Tests in nodes can be of the form $x_j > constant$
- Divides the space into axis-aligned rectangles

- Tests in nodes can be of the form $x_i > constant$
- Divides the space into axis-aligned rectangles



- Tests in nodes can be of the form $x_i > constant$
- Divides the space into axis-aligned rectangles



- Tests in nodes can be of the form $x_i > constant$
- Divides the space into axis-aligned rectangles



 Consider a split between each point in each dimension



6.034 - Spring 03 • 53

 Consider a split between each point in each dimension



6.034 - Spring 03 • 54

 Consider a split between each point in each dimension



 Choose split that minimizes average entropy of child nodes





6.034 - Spring 03 • 57



6.034 - Spring 03 • 58



6.034 - Spring 03 • 59



6.034 - Spring 03 • 60



6.034 - Spring 03 • 61



6.034 - Spring 03 • 62





Heart Disease

• Best performance (.77) slightly worse than nearest neighbor (.81)



6.034 - Spring 03 • 65

Heart Disease



thal = 1: normal exercise thallium scintigraphy test

6.034 - Spring 03 • 66



thal = 1: normal exercise thallium scintigraphy test ca = 0: no vessels colored by fluoroscopy



thal = 1: normal exercise thallium scintigraphy test ca = 0: no vessels colored by fluoroscopy










Auto MPG

 Performance (.91) essentially the same as nearest neighbor



6.034 - Spring 03 • 74

More than 22 MPG?



Bankruptcy Example



6.034 - Spring • 76

1-Nearest Neighbor Hypothesis



6.034 - Spring • 77

Decision Tree Hypothesis



Linear Hypothesis



6.034 - Spring • 79

Linearly Separable



Not Linearly Separable



Not Linearly Separable



Not Linearly Separable



Linear Hypothesis Class

• Equation of a hyperplane in the feature space

$$\mathbf{W} \cdot \mathbf{X} + b = 0$$
$$\sum_{j=1}^{n} w_j x_j + b = 0$$

• w, b are to be learned

Linear Hypothesis Class

• Equation of a hyperplane in the feature space

$$\mathbf{w} \cdot \mathbf{x} + b = 0$$
$$\sum_{j=1}^{n} w_j x_j + b = 0$$

• w, b are to be learned



Linear Hypothesis Class

• Equation of a hyperplane in the feature space

$$\mathbf{w} \cdot \mathbf{x} + b = 0$$

$$\sum_{j=1}^{n} w_j x_j + b = 0$$

• w, b are to be learned
• A useful trick: let x₀=1 and w₀=b

$$\overline{\mathbf{W}} \cdot \overline{\mathbf{X}} = \mathbf{0}$$
$$\sum_{j=0}^{n} W_j X_j = \mathbf{0}$$

• W,

Hyperplane: Geometry



Hyperplane: Geometry

 $\hat{\mathbf{w}} \cdot \mathbf{x} + b$

signed perpendicular distance of point **x** to hyperplane.



recall:
$$\mathbf{a} \cdot \mathbf{b} = \|\mathbf{a}\| \|\mathbf{b}\| \cos \theta$$

Hyperplane: Geometry



recall: $\mathbf{a} \cdot \mathbf{b} = \|\mathbf{a}\| \|\mathbf{b}\| \cos \theta$



Linear Classifier

Margin:

$$\gamma_i = y^i (\mathbf{w} \cdot \mathbf{x}^i + b) = y^i \overline{\mathbf{w}} \cdot \overline{\mathbf{x}}^i$$

proportional to perpendicular distance of point \mathbf{x}^{i} to hyperplane.

 $\gamma_i > 0$: point is correctly classified (sign of distance = y^i)

 $\gamma_i < 0$: point is incorrectly classified (sign of distance $\neq y^i$)



Perceptron Algorithm Rosenblatt, 1956

- Pick initial weight vector (including b), e.g. [0 ... 0]
- Repeat until all points correctly classified
 - Repeat for each point
 - -Calculate margin $(y^i \overline{wx})$ for point i
 - -If margin > 0, point is correctly classified
 - -Else change weights to increase margin; change in weight proportional to $y^i \overline{x}^i$

Perceptron Algorithm Rosenblatt, 1956

- Pick initial weight vector (including b), e.g. [0 ... 0]
- Repeat until all points correctly classified
 - Repeat for each point
 - -Calculate margin $(y^i \overline{wx})$ for point i
 - -If margin > 0, point is correctly classified
 - -Else change weights to increase margin; change in weight proportional to $y^i \overline{x}^i$
- Note that, if $y^i = 1$ if $x_j^i > 0$ then w_j increased (increases margin) if $x_j^i < 0$ then w_j decreased (increases margin)
- And, similarly for $y^i = -1$

Perceptron Algorithm Rosenblatt, 1956

- Pick initial weight vector (including b), e.g. [0 ... 0]
- Repeat until all points correctly classified
 - Repeat for each point
 - -Calculate margin $(y^{i}\overline{wx})$ for point i
 - -If margin > 0, point is correctly classified
 - -Else change weights to increase margin; change in weight proportional to $y^i \overline{x}^i$
- Note that, if $y^i = 1$ if $x_j^i > 0$ then w_j increased (increases margin) if $x_j^i < 0$ then w_j decreased (increases margin)
- And, similarly for $y^i = -1$
- Guaranteed to find separating hyperplane if one exists
- Otherwise, data are not linearly separable, loops forever

Perceptron Algorithm Bankruptcy Data



rate
$$\eta = 0.1$$

Initial Guess: w=[0.0 0.0 0.0]

Gradient Ascent

- Why pick $y^i \overline{\mathbf{x}}^i$ as increment to weights?
- To maximize scalar function of one variable f(w)
 - Pick initial w
 - Change w to w + η df/dw (η > 0, small)
 - until f stops changing $(df/dw \approx 0)$



Gradient Ascent/Descent

- To maximize $f(\mathbf{w})$ $\nabla_{\mathbf{w}} f = \left[\frac{\partial f}{\partial W_1}, \dots, \frac{\partial f}{\partial W_n}\right]$ • Pick initial \mathbf{w}
 - Change **w** to $\mathbf{w} + \eta \nabla_{\mathbf{w}} f$ ($\eta > 0$, small)
 - until f stops changing ($\nabla_{\mathbf{w}} f \approx 0$)
- Finds local maximum; global maximum if function is globally convex.

Gradient Ascent/Descent

- To maximize $f(\mathbf{w})$ $\nabla_{\mathbf{w}} f = \left[\frac{\partial f}{\partial W_1}, \dots, \frac{\partial f}{\partial W_n}\right]$ • Pick initial \mathbf{w}
 - Change **w** to $\mathbf{w} + \eta \nabla_{\mathbf{w}} f$ ($\eta > 0$, small)
 - until f stops changing $(\nabla_{\mathbf{w}} \mathbf{f} \approx \mathbf{0})$
- Finds local maximum; global maximum if function is globally convex
- Rate (η) has to be chosen carefully.
 - Too small slow convergence
 - Too big oscillation



Perceptron Training via Gradient Descent

• Maximize sum of margins of misclassified points

$$f(\mathbf{w}) = \sum_{i \text{ misclassified}} y^{i} \overline{\mathbf{w}} \overline{\mathbf{x}}^{i}$$

$$\nabla_{\mathbf{w}} f = \sum_{i \text{ misclassified}} y^i \overline{\mathbf{X}}^i$$

Perceptron Training via Gradient Descent

• Maximize sum of margins of misclassified points

$$f(\mathbf{w}) = \sum_{i \text{ misclassified}} y^{i} \overline{\mathbf{w}} \overline{\mathbf{x}}^{i}$$



- Off-line training: Compute gradient as sum over all training points.
- On-line training: Approximate gradient by one of the terms in the sum: yⁱxⁱ

Perceptron Algorithm Bankruptcy Data



Perceptron Algorithm Bankruptcy Data



Dual Form

Assume initial weights are 0; rate= $\eta > 0$



Dual Form

Assume initial weights are 0; rate= $\eta > 0$



$$h(\mathbf{x}) = sign(\overline{\mathbf{w}} \cdot \overline{\mathbf{x}}) = sign(\sum_{i=1}^{m} \alpha_{i} y^{i} \overline{\mathbf{x}}^{i} \cdot \overline{\mathbf{x}})$$

Perceptron Training Dual Form

- $\alpha = 0$
- Repeat until all points correctly classified
 - Repeat for each point i

-Calculate margin
$$\sum_{j=1}^{m} \alpha_{j} \mathbf{y}^{j} \overline{\mathbf{x}}^{j} \cdot \overline{\mathbf{x}}^{i}$$

-If margin > 0, point is correctly classified

– Else increment α_i

• Return
$$\overline{\mathbf{W}} = \sum_{j=1}^{m} \alpha_j \mathbf{y}^j \overline{\mathbf{X}}^j$$

 \bullet If data is not linearly separable, the α_{i} grow without bound