# CS325 Artificial Intelligence
## Chs. 18 & 4 – Supervised Machine Learning (cont)
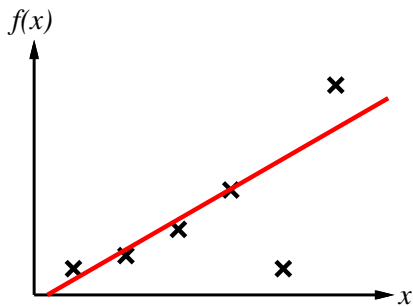
Cengiz Günay

Spring 2013

Let's start with the linear case. . .
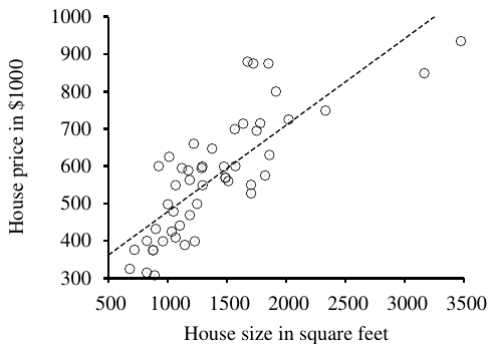
# Linear Regression

$$\text{price} = f(\text{size})$$

# Linear Regression



House price in \$1000 (y-axis), House size in square feet (x-axis)

$$\text{price} = f(\text{size})$$

$$? = f(3000)$$

$$y = f(x) = w_1 x + w_0$$

| $x$ | $y$ |
|----:|----:|
| 2 | 7 |
| $-1$ | $-2$ |
| 5 | 16 |
| $-3$ | $-8$ |

$w_0 = ?$
$w_1 = ?$

$$y = f(x) = w_1 x + w_0$$

| $x$ | $y$ |
|---|---|
| 2 | 7 |
| $-1$ | $-2$ |
| 5 | 16 |
| $-3$ | $-8$ |

$w_0 = 1$
$w_1 = 2$

$$y = f(x) = w_1 x + w_0$$

$$\mathrm{Loss}(f) = \sum_j \left(y_j - f(x_j)\right)^2 = \sum_j \left(y_j - (w_1 x_j + w_0)\right)^2$$

# Linear Regression—Defining a Loss Function

$$y = f(x) = w_1 x + w_0$$

$$\text{Loss}(f) = \sum_j (y_j - f(x_j))^2 = \sum_j (y_j - (w_1 x_j + w_0))^2$$

Minimum is where the derivative is zero:

$$\frac{\partial}{\partial w_0} \sum_j (y_j - (w_1 x_j + w_0))^2 = 0, \quad \frac{\partial}{\partial w_1} \sum_j (y_j - (w_1 x_j + w_0))^2 = 0$$

$$y = f(x) = w_1 x + w_0$$

$$\mathrm{Loss}(f) = \sum_j (y_j - f(x_j))^2 = \sum_j (y_j - (w_1 x_j + w_0))^2$$

Minimum is where the derivative is zero:

$$\frac{\partial}{\partial w_0} \sum_j (y_j - (w_1 x_j + w_0))^2 = 0, \quad \frac{\partial}{\partial w_1} \sum_j (y_j - (w_1 x_j + w_0))^2 = 0$$
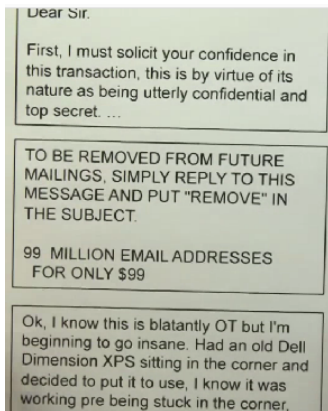
Solution is:

$$w_1 = \frac{N(\sum x_j y_j) - (\sum x_j)(\sum y_j)}{N(\sum x_j^2) - (\sum y_j)^2}, \quad w_0 = \left( \sum y_j - w_1(\sum x_j) \right) / N$$

- Everybody loves spam!



Dear Sir,

First, I must solicit your confidence in this transaction, this is by virtue of its nature as being utterly confidential and top secret. ...

TO BE REMOVED FROM FUTURE MAILINGS, SIMPLY REPLY TO THIS MESSAGE AND PUT "REMOVE" IN THE SUBJECT.

99 MILLION EMAIL ADDRESSES FOR ONLY $99

Ok, I know this is blatantly OT but I'm beginning to go insane. Had an old Dell Dimension XPS sitting in the corner and decided to put it to use, I know it was working pre being stuck in the corner,

Let's guess: $P(S) = \pi$

$$P(y_i) = \begin{cases} \pi & \text{if } y_i = S \\ 1 - \pi & \text{if } y_i = H \end{cases}$$

Let's guess: $P(S) = \pi$

$$P(y_i) = \begin{cases} \pi & \text{if } y_i = S \\ 1 - \pi & \text{if } y_i = H \end{cases}$$

Joint probability

$$\begin{aligned} P(\text{data}) &= \pi^{\text{count}(S)} \times (1 - \pi)^{\text{count}(H)} \\ &= \pi^3 \times (1 - \pi)^5 \end{aligned}$$

Joint probability

$$
\begin{aligned}
P(\mathrm{data}) &= \pi^{\mathrm{count}(S)} \times (1 - \pi)^{\mathrm{count}(H)} \\
&= \pi^3 \times (1 - \pi)^5
\end{aligned}
$$

Take log of both sides

$$
\log P(\mathrm{data}) = 3 \log \pi + 5 \log(1 - \pi)
$$

Joint probability

$$
\begin{aligned}
P(\mathrm{data}) &= \pi^{\mathrm{count}(S)} \times (1-\pi)^{\mathrm{count}(H)} \\
&= \pi^3 \times (1-\pi)^5
\end{aligned}
$$

Take log of both sides

$$
\log P(\mathrm{data}) = 3 \log \pi + 5 \log(1-\pi)
$$

Find max by zero derivative

$$
\frac{\nabla P(\mathrm{data})}{\nabla \pi} = 0 = \frac{3}{\pi} - \frac{5}{1-\pi}
$$

$$
\pi = 3/8
$$

# Bag of Words Representation

BAG OF WORDS

HELLO I WILL SAY HELLO

HELLO    I    WILL    SAY    } DICTIONARY
2        1     1       1

SPAM
OFFER IS SECRET
CLICK SECRET LINK
SECRET SPORTS LINK

HAM
PLAY SPORTS TODAY
WENT PLAY SPORTS
SECRET SPORTS EVENT
SPORT IS TODAY
SPORT COSTS MONEY

QUIZ    SIZE OF VOCABULARY =

SPAM

OFFER is SECRET
CLICK SECRET LINK
SECRET SPORTS LINK

HAM

PLAY SPORTS TODAY
WENT PLAY SPORTS
SECRET SPORTS EVENT
SPORT is TODAY
SPORT COSTS MONEY

Quiz ML-SOLUTIONS FOR

$P(\text{"SECRET"} \mid \text{SPAM}) = \boxed{\phantom{xxx}}$

$P(\text{"SECRET"} \mid \text{HAM}) = \boxed{\phantom{xxx}}$

SPAM

OFFER is SECRET
CLICK SECRET LINK
SECRET SPORTS LINK

HAM

PLAY SPORTS TODAY
WENT PLAY SPORTS
SECRET SPORTS EVENT
SPORTS is TODAY
SPORTS COSTS MONEY

Quiz  MESSAGE  M = "SPORTS"

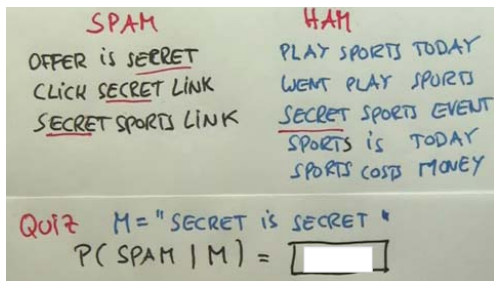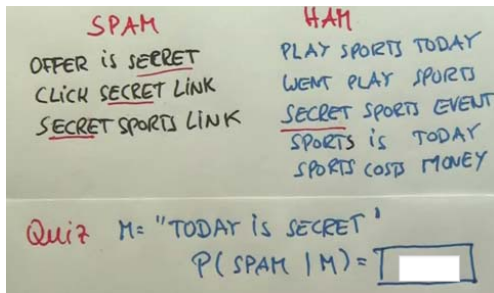$P(SPAM \mid M) = \boxed{\phantom{XXX}}$

$$P(S|M) = \alpha P(M|S)P(S)$$

$$P(S|M) = \alpha P(M|S)P(S)$$

$$P(S|M) = \alpha P(M|S)P(S) = \alpha P(M_1, M_2, M_3|S)P(S)$$

SPAM
OFFER is SECRET
CLiCK SECRET LINK
SECRET SPORTS LINK

HAM
PLAY SPORTS TODAY
WENT PLAY SPORTS
SECRET SPORTS EVENT
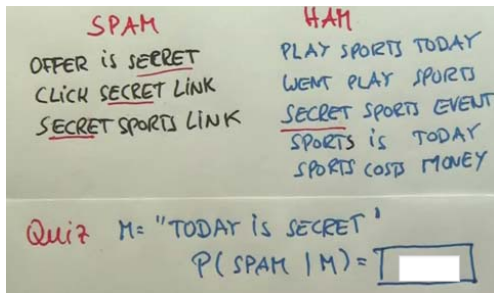SPORTS is TODAY
SPORTS COSTS MONEY

Quiz $M =$ "TODAY is SECRET"

$P(SPAM \mid M) = $ ☐

$$P(S|M) = 0$$

# Problems?



$$P(S|M) = 0$$

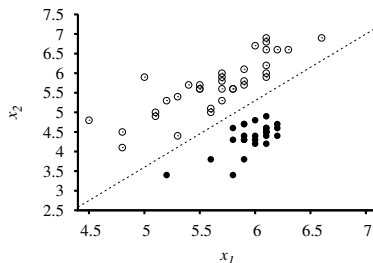- Need Laplace Smoothing (check the videos)

So far we calculated directly from data:

- **Linear regression coefficients** through explicit solution
- **Bayes net parameters** through maximal likelihood

So far we calculated directly from data:

- **Linear regression coefficients** through explicit solution
- **Bayes net parameters** through maximal likelihood

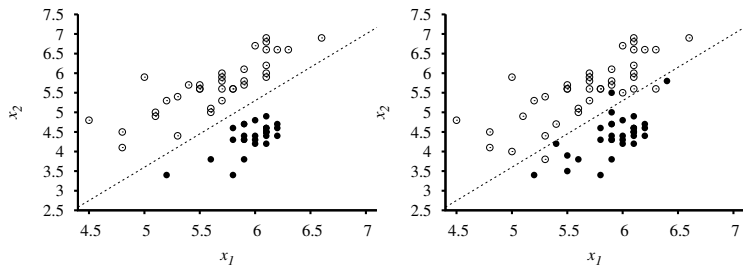But cannot solve every problem with these. Classification solution is not unique:
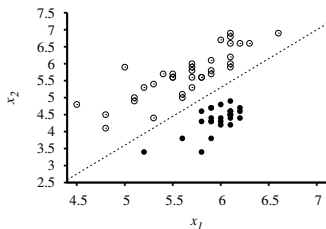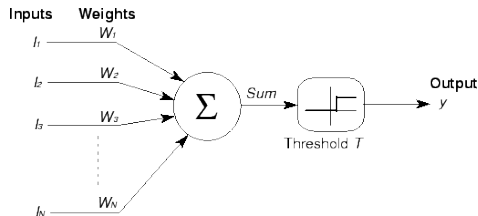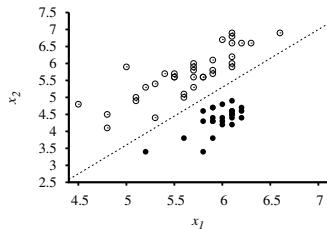
So far we calculated directly from data:

- **Linear regression coefficients** through explicit solution
- **Bayes net parameters** through maximal likelihood

But cannot solve every problem with these. Classification solution is not unique:

$$\text{sum} = \sum_{i=1}^{N} I_i W_i$$

$$y = \begin{cases} 1, & \text{if sum} \geq T \\ 0, & \text{if sum} < T \end{cases}$$
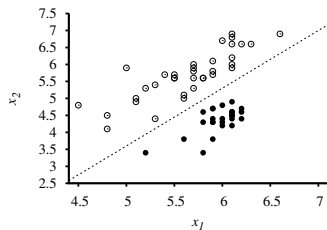
Line equation in 2D:

$$x_2 = a\,x_1 + b$$

Line equation in 2D:

$$x_2 = a\,x_1 + b$$
$$-b = a\,x_1 - x_2$$

Line equation in 2D:

$$x_2 = a\,x_1 + b$$
$$-b = a\,x_1 - x_2$$



The perceptron boundary:

$$T = w_1 x_1 + w_2 x_2$$

Line equation in 2D:

$$x_2 = a x_1 + b$$
$$-b = a x_1 - x_2$$



The perceptron boundary:

$$T = w_1 x_1 + w_2 x_2$$

$$y = \begin{cases} 1, & \text{if sum} \geq T \\ 0, & \text{if sum} < T \end{cases}$$

# Perceptron, 2D Case

Line equation in 2D:



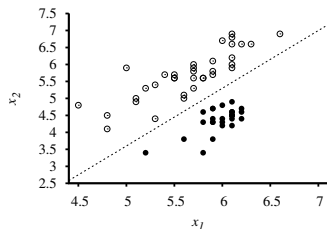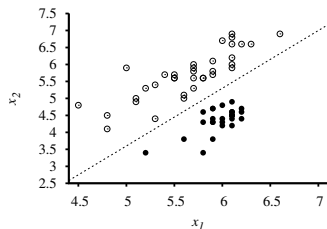$$x_2 = a x_1 + b$$
$$-b = a x_1 - x_2$$

The perceptron boundary:

$$T = w_1 x_1 + w_2 x_2$$

$$y = \begin{cases} 1, & \text{if sum} \geq T \\ 0, & \text{if sum} < T \end{cases}$$

How to learn it?

## Use the Loss Function, Perceptron



Perceptron:

$$y = f_{\mathbf{w}}(\mathbf{x})$$

Over all samples:

$$\text{Loss}(\mathbf{w}) = \sum_i \left(y_i - f_{\mathbf{w}}(\mathbf{x}_i)\right)^2$$

Perceptron:

$$y = f_{\mathbf{w}}(\mathbf{x})$$

Over all samples:

$$\text{Loss}(\mathbf{w}) = \sum_i (y_i - f_{\mathbf{w}}(\mathbf{x}_i))^2$$

Trying to find

$$\arg\min_{\mathbf{w}} \text{Loss}(\mathbf{w})$$

Perceptron:

$$y = f_{\mathbf{w}}(\mathbf{x})$$

Over all samples:

$$\text{Loss}(\mathbf{w}) = \sum_i (y_i - f_{\mathbf{w}}(\mathbf{x}_i))^2$$

Trying to find

$$\arg \min_{\mathbf{w}} \text{Loss}(\mathbf{w})$$

An incremental rule:

$$w_j \leftarrow w_j - \alpha \frac{\partial}{\partial w_j} \text{Loss}(\mathbf{w})$$

Perceptron:



$$y = f_{\mathbf{w}}(\mathbf{x})$$

Over all samples:

$$\text{Loss}(\mathbf{w}) = \sum_i \left(y_i - f_{\mathbf{w}}(\mathbf{x}_i)\right)^2$$
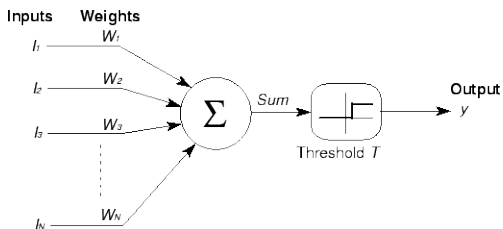
Trying to find

$$\arg \min_{\mathbf{w}} \text{Loss}(\mathbf{w})$$

An incremental rule:

$$w_j \leftarrow w_j - \alpha \frac{\partial}{\partial w_j} \text{Loss}(\mathbf{w})$$

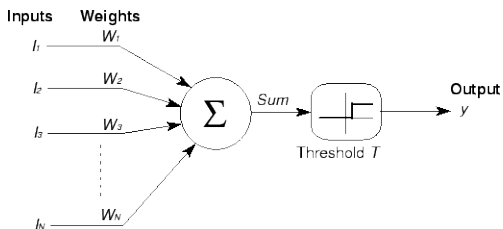$$w_j \leftarrow w_j + \alpha(y - f_{\mathbf{w}}(\mathbf{x})) \times x_j$$

$$w_j \leftarrow w_j + \alpha(y - f_{\mathbf{w}}(\mathbf{x})) \times x_j$$

|   |         | $y = 0$ | $y = 1$ |
|---|---------|---------|---------|
|   |         | Tom     | Jerry   |
| x | Trucks  | 1       | 0       |
|   | Sedans  | 0       | 1       |
|   | Hybrids | 0       | 1       |
|   | SUVs    | 1       | 0       |

$$w_j \leftarrow w_j + \alpha(y - f_{\mathbf{w}}(\mathbf{x})) \times x_j$$

|   |   | $y = 0$ | $y = 1$ |
|---|---|---------|---------|
|   |   | Tom | Jerry |
| x | Trucks | 1 | 0 |
|   | Sedans | 0 | 1 |
|   | Hybrids | 0 | 1 |
|   | SUVs | 1 | 0 |

**Start:** $\mathbf{w} = 0$, $\alpha = 1$, $T = 1$.
For $y_{\text{Tom}}$:

Inputs | Weights

$I_1$ — $W_1$

$I_2$ — $W_2$

$I_3$ — $W_3$

$I_N$ — $W_N$

$\Sigma$ — Sum — Threshold $T$ — Output $y$

$$w_j \leftarrow w_j + \alpha(y - f_{\mathbf{w}}(\mathbf{x})) \times x_j$$

|   |         | $y = 0$ | $y = 1$ |
|---|---------|---------|---------|
|   |         | Tom     | Jerry   |
| x | Trucks  | 1       | 0       |
|   | Sedans  | 0       | 1       |
|   | Hybrids | 0       | 1       |
|   | SUVs    | 1       | 0       |

**Start:** $\mathbf{w} = 0$, $\alpha = 1$, $T = 1$.

For $y_{\text{Tom}}$:

$$w_{\text{Trucks}} \leftarrow w_{\text{Trucks}} + (0 - 0) \times 1$$
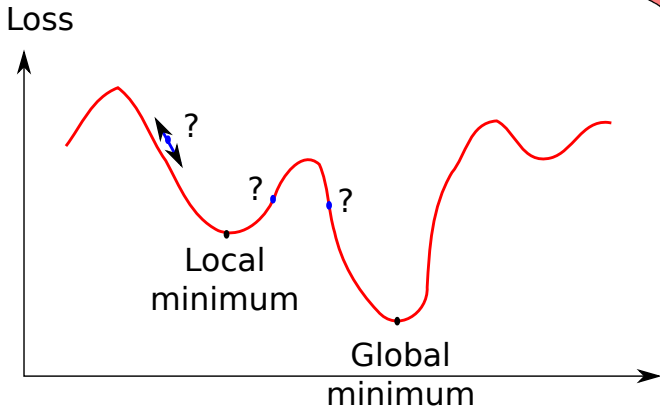
For $y_{Jerry}$:

$$w_j \leftarrow w_j + \alpha(y - f_{\mathbf{w}}(\mathbf{x})) \times x_j$$

|  |  | $y = 0$ | $y = 1$ |
|---|---|---|---|
|  |  | Tom | Jerry |
| | Trucks | 1 | 0 |
| $\mathbf{x}$ | Sedans | 0 | 1 |
| | Hybrids | 0 | 1 |
| | SUVs | 1 | 0 |

**Start:** $\mathbf{w} = 0$, $\alpha = 1$, $T = 1$.
For $y_{\text{Tom}}$:

$$w_{\text{Trucks}} \leftarrow w_{\text{Trucks}} + (0 - 0) \times 1$$

For $y_{\text{Jerry}}$:

$$w_{\text{Sedans}} \leftarrow w_{\text{Sedans}} + (1 - 0) \times 1$$

# Gradient Descent on the Loss Function

In general,

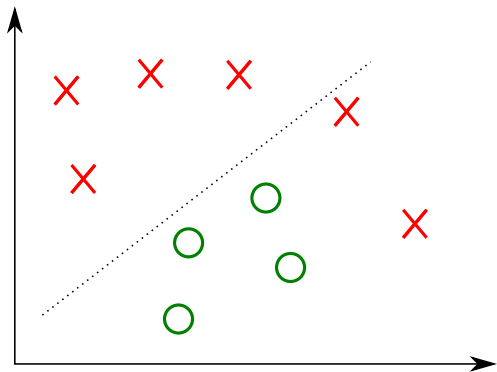$$w_j \leftarrow w_j + \alpha \frac{\partial}{\partial w_j} \text{Loss}(\mathbf{w})$$



Loss

?

?

?

Local
minimum

Global
minimum

In general,

$$w_j \leftarrow w_j + \alpha \frac{\partial}{\partial w_j} \mathrm{Loss}(\mathbf{w})$$



Loss



Local minimum

Global minimum

Adaptive $\alpha \Rightarrow$ **Simulated Annealing**

In general,

$$w_j \leftarrow w_j + \alpha \frac{\partial}{\partial w_j} \text{Loss}(\mathbf{w})$$



Loss

? 

? ?

Local
minimum

Global
minimum

Adaptive $\alpha$ $\Rightarrow$ **Simulated Annealing**

Major problem: **local minima**

# What If the Boundary is Non-linear?

# What If the Boundary is Non-linear?
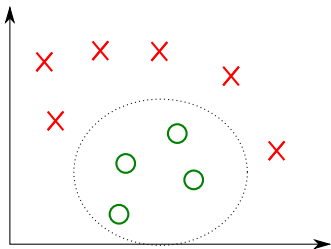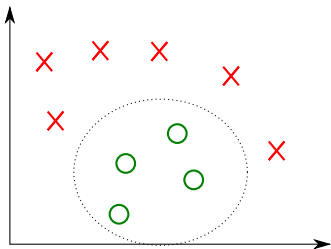
# What If the Boundary is Non-linear?
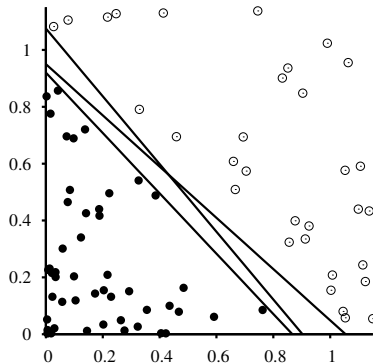


$\Rightarrow$ **Multi-Layer Perceptrons**

# Another Solution: Non-linear Kernels
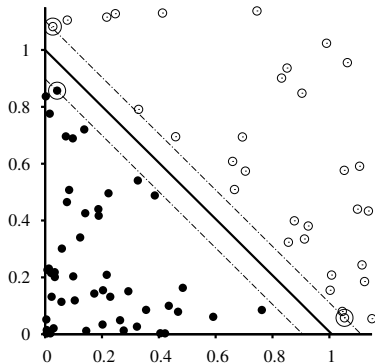


- Convert feature (input) space using non-linear kernel (e.g., radial distance)

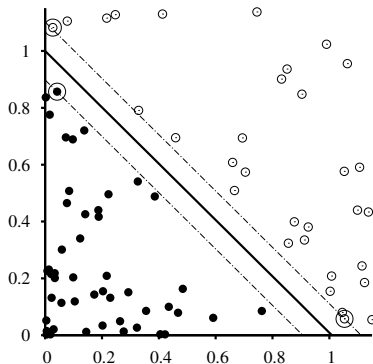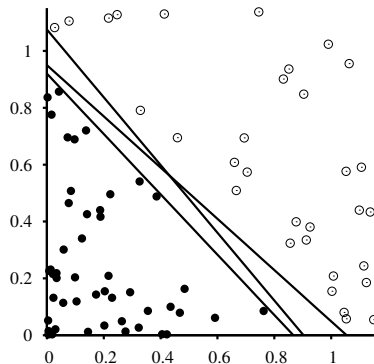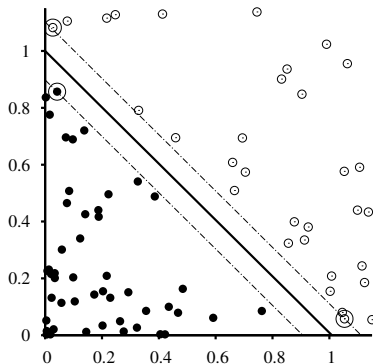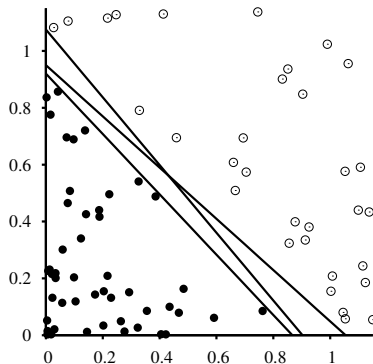- SVMs are guaranteed to find **optimal** solution $\Rightarrow$ Statistical Learning Theory

# Optimal Boundary? Enter Support Vector Machines



- SVMs are guaranteed to find **optimal** solution ⇒ Statistical Learning Theory
- **Kernel SVMs** are especially powerful because it can search in multi-dimensional kernel space

- Problem choosing model complexity
- Kernel type
- Structural complexity of MLP or SVM

- Problem choosing model complexity
- Kernel type
- Structural complexity of MLP or SVM

Solution, ask the data:

1. cross-validation
   Divide data into three sets: training, validate, test

- Problem choosing model complexity
- Kernel type
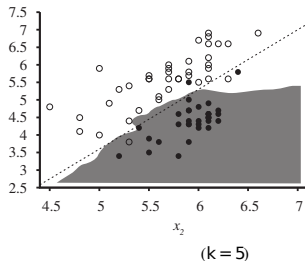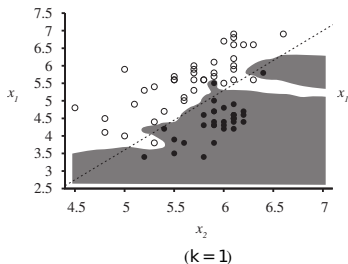- Structural complexity of MLP or SVM

Solution, ask the data:

1. cross-validation
   Divide data into three sets: training, validate, test

2. regularization
   Add complexity minimization term to Loss function

$$\text{Loss} = \sum (y_i - f(x_i))^2 + \beta \times \text{num params}$$

$k$-Nearest Neighbors algorithm:

- Keep all data points as lookup table
- Smoothing parameter, $k$



$(k = 1)$          $(k = 5)$

$k$-Nearest Neighbors algorithm:

- Keep all data points as lookup table
- Smoothing parameter, $k$



$(k = 1)$          $(k = 5)$

Problems?

$k$-Nearest Neighbors algorithm:

- Keep all data points as lookup table
- Smoothing parameter, $k$



Problems?

- Number of data points
- Number of features

Problems:

- No local minima, takes longer, must design problem well

# Summary of Supervised Machine Learning

- Can solve problems too complex for man-made algorithms
- Gets better with data (good for information age)
- Supervised learning with labels: regression and classification

## Summary of Supervised Machine Learning

- Can solve problems too complex for man-made algorithms
- Gets better with data (good for information age)
- Supervised learning with labels: regression and classification
- Linear regression and Bayes nets calculated from data directly
- Classification by minimizing Loss function iteratively

# Summary of Supervised Machine Learning

- Can solve problems too complex for man-made algorithms
- Gets better with data (good for information age)
- Supervised learning with labels: regression and classification
- Linear regression and Bayes nets calculated from data directly
- Classification by minimizing Loss function iteratively
- Local minima is a problem with gradient descent
- Non-linear problems can be solved with multiple boundaries or kernels
- Support vector machines find optimal solution faster

# Summary of Supervised Machine Learning

- Can solve problems too complex for man-made algorithms
- Gets better with data (good for information age)
- Supervised learning with labels: regression and classification
- Linear regression and Bayes nets calculated from data directly
- Classification by minimizing Loss function iteratively
- Local minima is a problem with gradient descent
- Non-linear problems can be solved with multiple boundaries or kernels
- Support vector machines find optimal solution faster
- Parameter complexity can be reduced with cross validation and regularization
- Non-parametric models good for low-dimensional problems
- Genetic algorithms have no local minima