\* A kernel will be valid if there exists a space such that

$$K(x_1, x_2) = \phi(x_1)^T \phi(x_2) = \sum_{n=1}^{N} \phi_n(x_1) \phi_n(x_2)$$

\* consider a quadratic kernel, with $D=2$,

$$K(x, z) = (x^T z)^2 = (x_1 z_1 + x_2 z_2)^2$$

$$= (x_1^2 z_1^2 + 2 x_1 z_1 x_2 z_2 + x_2^2 z_2^2)$$

$$x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

$$z = \begin{pmatrix} z_1 \\ z_2 \end{pmatrix}$$

$$x^T = (x_1, x_2)$$

$$x^T z = (x_1, x_2)\begin{pmatrix} z_1 \\ z_2 \end{pmatrix}$$

This can be expressed as an inner product of space where,

$$\phi(x) = x_1^2 + \sqrt{2} x_1 x_2 + x_2^2$$

this, gives, $\phi(z) = z_1^2 + \sqrt{2} z_1 z_2 + z_2^2$

$$\phi(x)\phi(z) = x_1^2 z_1^2 + 2 x_1 x_2 z_1 z_2 + x_2^2 z_2^2$$

$$K(x, z) = \phi(x)^T \phi(z)$$

\* A necessary and sufficient condition for a kernel function to be "valid" is that the gram matrix be positive and semidefinite for all choices of $\{x_m\}$.

A Gram matrix of $x$ is $x^T x$.

The linear vector $x$ is projected into a quadratic surface. If all the points in this surface are non-zero then our kernel is valid.

**\* Fisher's discriminant cost**

$$J(w) = Tr\{ [w \, S_W \, w^T]^{-1} \} \cdot (w \, S_B \, w^T)$$

— x — x — x — x — x — | Tue Oct 31 |

**Least square perceptron**    ← DOES not work

$$E(w) = \frac{1}{2} \sum_{n=1}^{N} (h_n - t_n)^2$$

$0 \quad \frac{1}{2} \cdots \frac{N-2}{2} \quad \frac{N}{2}$    N mistakes

$\frac{N+1}{2}$ mistake    →    N+1 values

min $E = 0$
max $E = \frac{N}{2}$

we cannot compute gradient, since function is
discrete and not continuous.

cost = No. of misclassified patterns
= 0 for NO mistake
$\frac{N}{2}$ for N mistakes

$cost = [0 \; \frac{1}{2} \; 2 \cdots \frac{N-2}{2} \; \frac{N}{2}]$    discrete set.



$t = sgn$

+1

0 ──────→ 3

−1

# Distance metrics

① Euclidean $\quad d(x,y) = \|x-y\|_2 = \sqrt{(x-y)^T (x-y)}$

② Hamming $\quad d(x,y) = $ # of different values in fixed length strings

③ Mahalanobis distance $\quad d(x,y) = \sqrt{(x-y)^T S^{-1} (x-y)}$

$S$ is sample cov. matrix

$S = I \rightarrow$ Euclidean

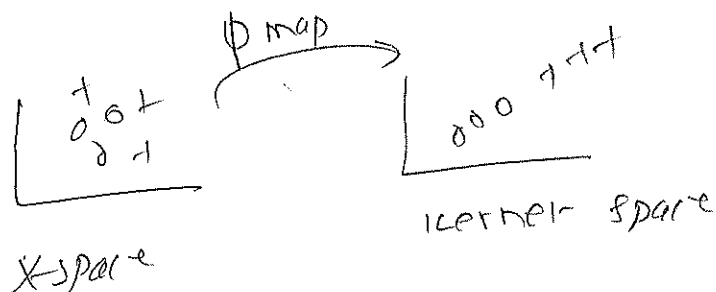$S = diag(\sigma_1^{-2}, \sigma_2^{-2}, ...) \rightarrow$ normalized Euclidean dist.

③ cosine similarity

$$d(x,y) = 1 - \cos(x,y) = 1 - \frac{x^T y}{\|x\| \|y\|}$$

④ Levenstein distance (edit distance)

min # of basis operations (del, insert, juxtapose) b/n two strings

$x = $ 'attens' $\quad y = $ 'hints' $\quad d(x,y) = 4$

$\phi$ map

$\begin{bmatrix} 1 & 0 & 1 \\ 0 & 0 & 1 \\ 0 & 1 \end{bmatrix}$ $\xrightarrow{\quad}$ $\begin{bmatrix} 0 & 0 & 0 & 1 & 1 & 1 \end{bmatrix}$

X-space $\qquad\qquad$ kernel space

## ① SVM with slack (soft margin SVM)

$$\min \quad J(\omega, b, \xi) = \frac{1}{2}\|\omega\|^2 + c\sum_{n=1}^{N}\xi_n$$

with constraint

$$t_n(\omega^T\phi(x_n) + b) \geq 1 - \xi_n$$

$$\xi_n \geq 0 \quad \text{and} \quad \sum_{i=1}^{N}\xi_i \leq Z \qquad \forall n \in \{1 \cdots N\}$$

i.e.

$$1 - \xi_n - t_n\omega^T\phi(x_n) - bt_n \leq 0 \quad \text{——①} \quad \text{(convex fn)}$$

$$-\xi_n \leq 0 \quad \text{——②} \quad \text{(convex fn constraint)}$$

$$\forall n = 1, 2, \cdots, N$$

## primal Lagrangian

$$L_p = L(\omega, b, \xi, \alpha, r) = \frac{1}{2}\|\omega\|^2 + c\sum_{n=1}^{N}\xi_n$$

$$+ \sum_{n=1}^{N}\alpha_n(1 - \xi_n - t_n\omega^T\phi_n - bt_n)$$

$$- \sum_{n=1}^{N}r_n\xi_n$$

## dual Lagrangian

$$L_g(\alpha, r) = \inf_{\omega, b, \xi} L_p(\omega, b, \xi, \alpha, r)$$

Node I

$$0 = \frac{\partial L_p(\omega, b, \xi, \alpha, r)}{\partial \omega} = \omega - \sum_n \alpha_n t_n \phi(x_n) \Rightarrow \boxed{\omega = \sum_n \alpha_n t_n \phi_n}$$

$$0 = \frac{\partial L_p}{\partial b} = -\sum_n \alpha_n t_n \Rightarrow \boxed{\sum_n \alpha_n t_n = 0}$$

NO summation!

$$0 = \frac{\partial L_p}{\partial \xi_n} = c - \alpha_n - r_n \Rightarrow \boxed{c = \alpha_n + r_n} \quad \forall n \in \{1, 2, \cdots N\}$$

Then,

$$L_D(\xi, r) = \frac{1}{2}\|w\|^2 + c\sum_n \xi_n + \sum_n \alpha_n(1 - \xi_n - t_n w^T \phi_n - t_n b) - \sum_n r_n \xi_n$$

$$L_D = \sum_n \alpha_n - \frac{1}{2}\sum_{m,n=1}^{N} \alpha_m \alpha_n t_m t_n k(x_m, x_n)$$

Then the optimization problem in dual spareis,

$$\underset{\alpha}{\text{maximize}} \quad L_D(\alpha), = \sum_n \alpha_n - \frac{1}{2}\sum_{m,n} \alpha_m \alpha_n t_m t_n k(x_m, x_n)$$
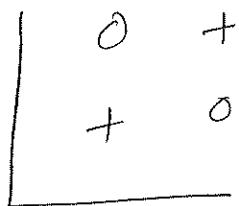
with constraints

$$0 \le \alpha_n \le c \quad \forall n = 1, 2, \ldots, N$$

$$\sum_{n=1}^{N} \alpha_n t_n = 0$$

note: $c\sum_n \xi_n - \sum_n \alpha_n \xi_n - \sum_n r_n \xi_n = 0$

since, $c\sum_n \xi_n = \sum_n \alpha_n \xi_n + \sum_n r_n \xi_n$

$$= \sum_n (\alpha_n + r_n)\xi_n$$

(1)

| | |
|---|---|
| 0 | + |
| + | 0 |

- SVM (quad kernel) can achieve zero training error
- Logis Regr, 3-NN cannot

(2) If examples are iid, increase training examples → may increase train error
↳ but decrease test error

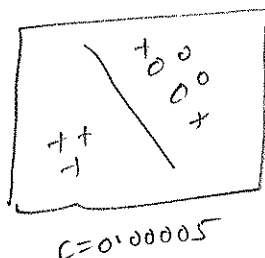(3) SVM effect of $c$

$$\min \; J(w,b) = \frac{1}{2}\|w\|^2 + c \sum_{n=1}^{N} \xi_n$$

s.t. $\quad w^T \phi(x_n) + b \geq 1 - \xi_n \quad \forall \; n \in \{0, N\}$

will not change
↳ decision bond.



$C = 10,000$

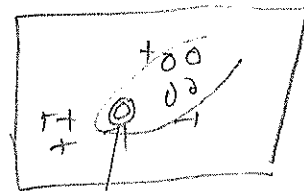$C = 0.00005$

$C = 1$

adding this change
dec. boun. drastically

between $c \gg 1$ and $c \approx 0$ choose $c \approx 0$ because it maximizes the margin between dominant cloud of points and we can not depend on any few data points which can be noise.

(4) Bias variance Tradeoff

| | Bias | var |
|---|---|---|
| linear regr | high | low |
| $d=2$ poly | low | low |
| $d=10$ poly | low | high |

# python

① $X = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}$     append ones column to data X.

$X = np.array([\ [1,2,3], [4,5,6]\ ]) = np.arange(1,7).reshape(2,3)$

$ones = np.ones(X.shape[0]).reshape(-1,1)$

$X1 = np.append(ones, X, axis=1).astype(np.int)$
   concatenate

$X1 = np.c\_[np.ones(X.shape[0]).reshape(-1,1), X]$

$X1 = np.column\_stack((\ "\ ))$

$X1 = np.c\_[np.ones(X.shape[0])[:, np.newaxis], X]$
                                        .reshape(-1,1)

$= np.c\_[np.atleast\_2d(np.ones(X.shape[0]).T, X]$

$= np.c\_[np.expand\_dims(np.ones(X.shape[0], axis=1), X]$

⑤ Given $\phi(x) = [1, x_1, x_2, x_1 x_2]$

  find the kernel $K(x, x')$.

  $\Rightarrow K(x, x') = 1 + x_1 x_1' + x_2 x_2' + x_1 x_2 x_1' x_2'$

⑥ L1 VS L2 loss

  ⓐ False: L2 is more robust to outlier than L1
  gradient of L2 loss can grow without bounds, but
  gradient of L1 loss is bounded, hence
  influence of outlier is limited.

  (Note: L2 gives more values to misclassification
    than L1)

  ⓑ L1 gives sparse solution & used in feature selection.

  ⓒ Logistic loss is better than L2 loss in classification task.

① SVM small C,

  For linearly separable data, small C can affect
  training accuracy.
      A small C can allow large slacks, thus, the
  resulting classifier will have smaller $w^2$
  and can have non-zero training error.

# Q1 ① solve the SVM problem without slack using Lagrange multiplier method

→ The optimization problem is

$$\text{minimize} \quad J(w,b) = \frac{1}{2}\|w\|^2$$

$$\text{s.t.} \quad t_n(w^T\phi(x_n)+b) \geq 1 \quad \forall n \in \{1,\dots,N\}$$

$$1 \leq t_n(w^T\phi(x_n)+b)$$

$$1 \leq t_n w^T\phi(x_n) + t_n b$$

$$\underline{1 - t_n w^T\phi(x_n) - t_n b \leq 0} \quad (\text{convex fn constraint})$$

$$\text{compare} \quad f_i(x) \leq 0 \quad \text{then} \quad i = 1,\dots,m$$

primal Lagrangian,

$$L_p(w,b,\alpha) = \frac{1}{2}\|w\|^2 + \sum_{n=1}^{N}\alpha_n(1 - t_n w^T\phi_n - t_n b)$$

where $\alpha_n \geq 0$ are Lagrange multipliers.

Dual Lagrangian,

$$L_D(\alpha) = \inf_{w,b} L_p(w,b,\alpha)$$

first find the infimum of $L_p$ w.r.t $w,b$:

$$\frac{\partial}{\partial w} L_p = 0 = w + \sum_n(-)t_n\phi_n \quad \Rightarrow \quad \boxed{w = \sum_n \alpha_n t_n \phi_n} \quad \text{---①}$$

$$\frac{\partial}{\partial b} L_p = 0 = \sum_n(-)\alpha_n t_n \quad \Rightarrow \quad \boxed{\sum_n \alpha_n t_n = 0} \quad \text{---②}$$

LOOK LHS

Then Dual Lagrangian is,

$$L_D(\alpha) = \frac{1}{2} \sum_{n,m} \alpha_m \alpha_n t_m t_n \phi_m \phi_n + \sum_n \alpha_n - \sum_n \alpha_n t_n \phi_n \cdot \sum_m \alpha_m t_m \phi_m$$

$$- \sum_n \alpha_n t_n \overset{0}{\cancel{b}}$$

$$\boxed{L_D(\alpha) = \sum_n \alpha_n - \frac{1}{2} \sum_{n,m} \alpha_m \alpha_n t_m t_n \, K(x_n, x_m)}$$

$$K(x,y) = \vec{\phi_{(x)}} \cdot \vec{\phi_{(y)}} = \phi_{(x)}^T \phi_{(y)}$$

Then the optimization problem in Dual space is,

maximize $L_D(\alpha) = \sum_n \alpha_n - \frac{1}{2} \sum_n \sum_m \alpha_m \alpha_n t_m t_n \, K(x_m, x_n)$

s.t. $\alpha_n \geq 0 \quad \forall n \in \{1, \dots N\}$

$$\sum_{n=1}^{N} \alpha_n t_n = 0$$

Look here ↑

Now, KKT conditions are,

note!
we can write
$y_{(x_n)} = w^T \phi_{(x_n)} + b$

① primal constraints $\quad 1 - t_n(w^T \phi_{(x_n)} + b) \leq 0 \quad$ (convex constraint equation)
$$t_n(w^T \phi_{(x_n)} + b) \geq 1$$

② dual constraint $\quad \alpha_n \geq 0 \quad \forall n \in \{1, \dots N\}$

③ complementary slackness $\quad \alpha_n \{1 - t_n w^T \phi_{(x_n)} - t_n b\} = 0$

for any data point, either $\alpha_n = 0$

$$1 - t_n w^T \phi_{(x_n)} - t_n b = 0$$
↑
these $\alpha_n$ are called
support vectors

here $n = 1, 2, \ldots, N$

Suppose, out of N samples, there are m support vectors then N-m samples will have Lagrange parameter $\alpha = 0$ and m examples will have non-zero Lagrange parameters.

$$1 - t_m \, w^T \phi(x_m) - t_m b = 0$$

or, $\quad 1 - t_m \, \phi(x_m) \sum_n \alpha_n t_n \phi(x_n) - t_m b = 0$

or, $\quad t_m b = 1 - t_m \, \phi(x_m) \sum_n \alpha_n t_n \phi(x_n)$

$$= 1 - t_m \sum_n \alpha_n t_n \phi(x_n) \, \phi(x_m)$$

$$b \, t_m = 1 - t_m \sum_n \alpha_n t_n \phi(x_n) \, \phi(x_m)$$

$$b = \frac{1}{t_m} - \sum_n \alpha_n t_n \, \phi(x_n) \phi(x_m) = t_m - \sum_n \alpha_n t_n \phi(x_n) \cdot \phi(x_m)$$

$$\boxed{b = t_m - w \cdot \phi(x_m)}$$

$$\boxed{b = t_m - \sum_n \alpha_n t_n \, k(x_n, x_m)} \qquad \because \frac{1}{t_m} = t_m$$

$$1 = \frac{1}{1} \quad \text{and} \quad -1 = \frac{1}{-1}$$

this is true for all the m examples which have non-zero Lagrange parameter $\alpha$.

For numerical stability we choose value of b as the mean of all b-values, then,

$$\boxed{b = \frac{1}{|S|} \sum_{m \in S} \left[ t_m - \sum_{n \in S} \alpha_n t_n \, k(x_n, x_m) \right]}$$

where S is the subset of all the examples where Lagrange parameter $\alpha$ is non-zero.

$$S \subseteq D$$

$$S = \{ n \mid 1 - t_n \, w^T \phi(x_n) - t_n b = 0 \}$$

Then, Linear discriminant function is,

$$y(x) = w^T \phi(x) + b = \sum_{n \in S} \alpha_m t_m \, k(x, x_n) + \frac{1}{|S|} \sum_{m \in S} \left[ t_m - \sum_{n \in S} \alpha_n t_n k(x_n, x_m) \right]$$

$$\underline{KNN} = \text{memory-based (no model to fit)}$$

① $\begin{cases} \text{find k nearest examples } x_1, x_2 \cdots x_k \text{ from test } x_i \\ y(x) = \underset{t \in T}{\arg\max} \; \overset{k}{\underset{i=1}{\sum}} \delta_t(t_i) \end{cases}$

$w_i$ gives distance-weighted k-NN

$w_i = \frac{1}{\|x - x_i\|^2}$

② Mahalanobis dB

$$d(x,y) = \sqrt{(x-y)^T S^{-1} (x-y)}$$

sample covariance matrix
if $S = I$ = Euclidean distance
$S = \text{diag}(\sigma_1^{-2}, \sigma_2^{-2}, \cdots, \sigma_k^{-2})$
normalized Euclidean

③ cosine similarity

$$d(x,y) = 1 - \frac{x^T y}{\|x\| \|y\|}$$

④ kernel-based distance weighted NN
→ binary classification $T \in \{+1,-1\}$
→ $y(x) = \text{sign}\left( \overset{N}{\underset{i=1}{\sum}} k(x, x_i) \cdot t_i \right)$

# Wrapper method

① greedy
Forward selection
F = all features
S = subset of features
start  S = { } empty
for each feature f in F-S
find best performing feature f* and add to S.

Repeat until performance does not increse
or performance good enough


② Recursive
Backward Elimination

F: (1,2, ... K) is set of features
S = [ ] ranked set of features

Repeat until F-S is empty
train w using linear svm and F-S
find feature f with minimum |w|
append f to S

Return S

④ distance-weighted KNN for regression

1. find $k$ nearest points $x_{i1}, x_{i2} \dots x_{ik}$

2. $y(x) = \dfrac{\sum\limits_{i=1}^{k} w_i t_i}{\sum\limits_{i} w_i}$    where   $w_i = \dfrac{1}{\|x - x_i\|^2}$

   $k = N \Rightarrow$ Shepard method

$y(x) = \sum\limits_{i=1}^{N} K(x, x_i)\, t_i \;\Big/\; \sum\limits_{i=1}^{N} K(x, x_i)$   (kernel-based dist weighted)

⑤ Regression with KNN   (ti are values, not classes)

$y(x) = \dfrac{1}{k} \sum\limits_{i=1}^{k} t_i$

⑥ distance-weighted KNN (regression)

poional : $y(x) = \sum\limits_{i=1}^{k} w_i t_i \;\Big/\; \sum\limits_{i=1}^{k} w_i$     $w_i = \dfrac{1}{\|x - x_i\|^2}$

kernel based : $y(x) = \sum\limits_{i=1}^{N} K(x, x_i)\, t_i \;\Big/\; \sum\limits_{i=1}^{N} K(x, x_i)$   (ti are values not classes)

# 3 Approaches to RL



$\downarrow S$

direct use
ind. learning

$\boxed{\pi}$ $\xleftarrow{\text{argmax}}$ $\boxed{U}$ $\xleftarrow[\substack{\text{Bellman} \\ \text{eqn}}]{\text{solve}}$ $\boxed{T, R}$

$\downarrow a$

$\downarrow V$

$S \quad q$
$\downarrow \quad \downarrow$

$\downarrow \quad \downarrow$
$S' \quad r$

direct learning
indirect use

policy
search

value based
function
based

model-based

# Filter method of feature selection

## ① Mutual Information

$$MI(X,Y) = \sum_X \sum_Y p(x,y) \, \ln \frac{p(x,y)}{p(x)\, p(y)}$$

$= 0$ when $x, y$ indp
max when $x = y$

Let there are K examples with $l$ features.

$O_{ij}$ = observed value
for $x=i$, $y=j$

② 
$1 \cdots j \cdots l$

$\begin{matrix} 1 \\ i \\ K \end{matrix}$ $\quad O_{ij}$ $\quad \} N_{x=i}$

$E_{ij} = \dfrac{N_{x=i} \, N_{y=j}}{N}$

$\underbrace{\quad\quad}_{N_{y=j}}$

$$\chi^2 = \sum_{i=1}^{K} \sum_{j=1}^{l} \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

classification
@ probabilistic (Generative) models

{ Naive Bias
  Hidden Markov models

- inf, dec → separate          detection
- can use $p(x)$ for outlier or novelty
- need to model dependencies between features

Naive Bias → resilient to noise

- Text classification with NB → look HW

- multiple classes
  posterior prob of class $c_k$ given data $x$ is,

$$p(c_k | x) = \frac{p(x | c_k) \cdot p(c_k)}{\sum_j p(x | c_j) p(c_j)}$$


$p(x|y)$

$$= \frac{\exp(a_k(x))}{\sum_j \exp(a_j(x))}$$

normalized exponential

( softmax fn )

generative     where  $a_k(x) = \ln p(x | c_k) \cdot p(c_k)$

① SVM for ranking

optimization problem

minimize $J(w,b) = \frac{1}{2}\|w\|^2 + C\sum \xi_{k,i,j}$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \xi_{k,i,j} - z = 1 - \xi_{k,i,j}$

s.t. $w^T \phi(q_k, d_i) \geq w^T \phi(q_k, d_j) + 1 - \xi_{k,i,j}$

$\qquad \xi_{k,i,j} \geq 0$

( for a query $q_k$ we want document
$d_i$ be ranked higher than document $q_j$ )

① Add ones column
$= np.c_[np.ones(X.shape[0]).reshape(-1,1), X]$

$X = np.c_[np.ones(X.shape[0])[np.newaxis].T, X]$

② correct $= np.sum(y\_pred == y\_test)$

accuracy $= correct / len(y\_pred)$

① Gradient Descent ( GD or BGD)       GD with momentum

vanilla GD $\quad v^{\tau+1} = \eta \nabla J(w^\tau)$           $v^{\tau+1} = \gamma v^\tau + \eta \nabla J(w^\tau)$

$\quad\quad\quad w^{\tau+1} = w^\tau - v^{\tau+1}$             $w^{\tau+1} = w^\tau - v^{\tau+1}$

Batch Gradient Descent        (Lect 01, p.24)

$$J(w) = \frac{1}{2N} \sum_{n=1}^{N} (h_n - t_n)^2$$

$$w^{\tau+1} = w^\tau - \eta \, \nabla J(w^\tau)$$

$$= w^\tau - \eta \sum_{n=1}^{N} (h_n - t_n) \, x_n$$

$$= w^\tau - \text{learningRate} \, (h - t) \, @ \, X1$$

# confusion matrix

Actual class

predicted class cat

|  | cat | not cat |  |
|---|---|---|---|
| cat | TP | FP | $\tilde{P}$ |
| not cat | FN | TN | $\tilde{N}$ |

True

—————————

P    N

$00 = TN$
$01 = FP$
$10 = FN$
$11 = TP$

|  | 1 | 0 |  |
|---|---|---|---|
| 1 | TP | FP |  |
| 0 | FN | TN |  |

—————————

P    N

## actual

|  | 0 | 1 |  |
|---|---|---|---|
| 0 | TN | FN | $\tilde{N}$ |
| 1 | FP | TP | $\tilde{P}$ |

—————————

N    P

$\rightarrow$ precision $= \dfrac{TP}{TP+FP}$

$\text{recall} = \dfrac{TP}{P} = \dfrac{TP}{TP+FN}$

(hit rate)

$\dfrac{\text{predicted+ve}}{\text{actual +ve}}$

F1 score is the Harmonic mean of precision and recall.

$$F1 = \dfrac{2}{\dfrac{1}{\text{precision}} + \dfrac{1}{\text{recall}}}$$

$$= \dfrac{2\,\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

\* prove that the number of elements in x and y is also a kernel

$\Longrightarrow |K(x,y)| = |x \cap y|$ is a kernel.

— x — x — > — — x —————

**contd** (NULL)

WT $\phi(0) = \sum_{n=1}^{N} \alpha_n t_n K(x_n, x)$

strict
since

$s = n \quad r_n > 0$
$n \notin s \Rightarrow \alpha_n = 0$

$= \sum_{m \in S} \alpha_m t_m K(x_m, x)$

$+ \sum_{n \notin S} \alpha_m t_m K(x_m, x)$ but if $x_m \notin S$, then $\alpha_m = 0$
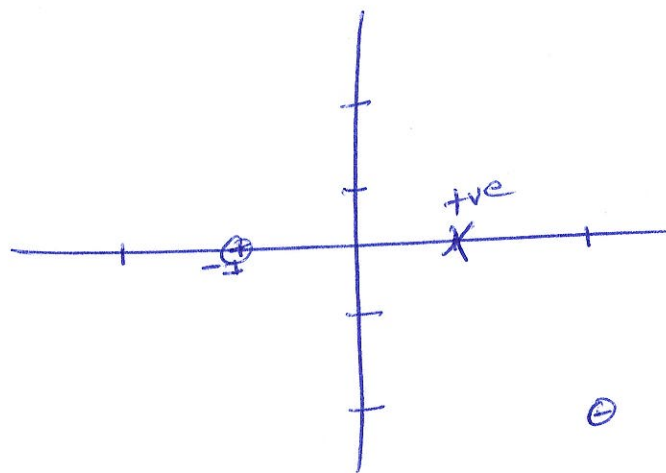
\* package libsum or sum lite,

input $\longrightarrow \{x_n, t_n)\}$ $1 \leq n \leq N$

output $\longrightarrow \{\alpha_m, t_m, x_m\}$ $m \in S$

model.txt $\rightarrow$ $\boxed{\begin{array}{c} b \\ \alpha_m t_m \cup x_m \end{array}}$

# perceptron update

$-1, 0 \quad -ve$
$2, -2 \quad -ve$ }

$1, 0 \quad \rightarrow +ve$



tve

$-1$

$\ominus$

final $w = 0\ 1\ 2$

let's say new point
$x = 2, -1\cdot01$
$x_{avg} = 1\ 2\ -1\cdot0$
$w \cdot x_{avg} = 0 + 2\ -2\cdot02$
$< 0$
$\therefore -ve$ label

## Augmented data

| | | | |
|---|---|---|---|
| 1 | -1 | 0 | -1 |
| 1 | 2 | -2 | -1 |
| 1 | 1 | 0 | 1 |

$x_{avg} \qquad t$

initial wt $\vec{w} = (0\ 0\ 0)$

note we can
also use
$\vec{w} \cdot \vec{x} + b$
$n = 0.9$ original unaugmented $x$
$b = 1$ etc

Test point    hypothesis correctly classified?    Updated weights correct example

| | Test point | | Updated weights |
|---|---|---|---|
| eg1 | $-: (1\ -1\ 0)$ | $wx = 0+0+0 \leq 0$   false | $w = w - x = (-1\ 1\ 0)$ |
| eg2 | $-: (1\ 2\ -2)$ | $wx = -1 + 2 + 0 < 0$ false | $w = w - x = (-2\ -1\ 2)$ |
| eg3 | $+: (1\ 1\ 0)$ | $wx = -2 - 1 + 0 \geq 0$ false | $w = w + x = (-1\ 0\ 2)$ |
| eg1 | $-: (1\ -1\ 0)$ | $wx = -1 + 0 + 0 < 0$ true | $w = w = (-1, 0\ 2)$ |
| eg2 | $-: (1\ 2\ -2)$ | $wx = -1 + 0 + (-4) < 0$ true | $w = w = (-1, 0, 2)$ |
| eg3 | $+: (1\ 1\ 0)$ | $wx = -1 + 0 + 0 \geq 0$ false | $w = w + x = (0\ 1\ 2)$ |
| eg1 | $-: (1 -1\ 0)$ | $wx = 0 + (-1) + 0 < 0$ true | $w = w = (0\ 1\ 2)$ } |
| eg2 | $-: (1\ 2 -2)$ | $wx = 0 + 2 + (-4) < 0$ true | $w = w = (0\ 1\ 2)$ } |
| eg3 | $+: (1\ 1\ 0)$ | $wx = 0 + 1 + 0 \geq 0$ true | $w = w = 0\ 1\ 2$ |

final weight

① constrained optimization
   ( Lagrange multipliers)

maximize $5 - (x_1-2)^2 - 2(x_2-1)^2$

s.t. $x_1 + 4x_2 = 3$

soln: If we ignore constraint we get $x_1 = 2, x_2 = 1$
      then $x_1 + 4x_2 = 2 + 4\cdot1 = 6$ is too large
      for the constraint.
      consider,
      $L = L(x_1, x_2, \lambda) = 5 - (x_1-2)^2 - 2(x_2-1)^2 + \lambda(3 - x_1 - 4x_2)$

LOOK $\lambda$: $\lambda = 0$    2, 1

$\lambda = 1$    3/2, 0

$\lambda = \frac{2}{3}\left(\frac{5}{3}, \frac{1}{3}\right)$    $\frac{5}{3} + 4\cdot\frac{1}{3} = \frac{5}{3} + \frac{4}{3} = \frac{9}{3} = 3$   ✓

formal soln) $\dfrac{\partial L}{\partial x_1} = -2(x_1-2) - \lambda = 0$

$= -2x_1 + 4 - \lambda = 0 \Rightarrow 2x_1 = 4 - \lambda$

$2x_1 = 4 - \lambda$
$= 4 - \frac{2}{3}$
$= \frac{4\cdot3 - 2}{3} = \frac{10}{3}$
$x_1 = \frac{5}{3}$

$\dfrac{\partial L}{\partial x_2} = -4(x_2-1) - 4\lambda = 0$

$= -4x_2 + 4 - 4\lambda = 0$

$\Rightarrow 4x_2 = 4 - 4\lambda$
$x_2 = 1 - \lambda$

$\dfrac{\partial L}{\partial \lambda} = 3 - x_1 - 4x_2 = 0$

$\Rightarrow 6 - 2x_1 - 8x_2 = 0$

$\Rightarrow 6 - 4 + \lambda - 8(1-\lambda) = 0$

$\Rightarrow 6 - 4 + \lambda - 8 + 8\lambda = 0$

$-6 + 9\lambda = 0$
$9\lambda = 6$
$\lambda = 2/3$

**OP1**

$$\min_{w,b} \tfrac{1}{2}\|w\|^2$$
$$s.t. \quad t_n\left(w^T\phi(x_n)+b\right)\geq 1$$

solution $= w_1^*, b_1^*$

**OP2**

$$\min_{w,b} \tfrac{1}{2}\|w\|^2$$
$$s.t. \quad t_n\left(w^T\phi(x_n)+b\right)\geq \eta$$

$\Downarrow$ OP3

$$\min_{w,b} \tfrac{1}{2}\left\|\tfrac{w}{\eta}\right\|^2$$
$$s.t. \quad t_n\left(\left(\tfrac{w}{\eta}\right)^T\phi(x_n)+\tfrac{b}{\eta}\right)\geq 1$$

OP3 has solution

$$w_2^* = w_1^*\,\eta$$
$$b_2^* = b_1^*\,\eta$$

$\Downarrow$ OP3 rename

$$\min_{w,b} \tfrac{1}{2}\|w'\|^2$$
$$s.t. \quad t_n\left(w'^T\phi(x_n)+b'\right)\geq 1$$

$$w' = w/\eta$$
$$b' = b/\eta$$

Now, decision hyperplanes are,

$$H1 = \left\{ x \mid w_1^{*T}\phi(x) + b_1^* = 0 \right\}$$

$$H2 = \left\{ x \mid w_2^{*T}\phi(x) + b_2^* = 0 \right\}$$
$$= \left\{ x \mid \eta w_1^{*T}\phi(x) + \eta b_1^* = 0 \right\}$$
$$= \left\{ x \mid w_1^{*T}\phi(x) + b_1^* = 0 \right\}$$

$$H2 = H1 \qquad q.e.d$$

③ kmcp  kernel multi-class perceptron

1  define $f(x) = \sum_{ij} \alpha_{ij} \left[ \phi(x_i, t_i)^T \phi(x, t) - \phi(x_i, c_j)^T \phi(x, t) \right]$

$\left( (\phi_f \omega)^T x = \sum_n \alpha_n t_n x_n^T x = \sum_n \alpha_n t_n K(x_n, x) \right)$

2  initialize dual params $\alpha_{ij} = 0$

3  for $i = 1 \cdots n$

4      $c_j = \underset{t \in T}{\operatorname{argmax}} \; f(x_i, t)$     ⎱ Repeat
         $(h_n = sgn(f(x)))$

5      if $c_j \neq t_i$ then     $h_n t_n$
          $\alpha_{ij} = \alpha_{ij} + 1$     $\alpha_n \alpha_{n+1}$

6

Testing:  $t^* = \underset{t \in T}{\operatorname{argmax}} \; f(x, t)$      $(h(x) = sgn(f(x)))$

② mcp    ( consept of kernel comes from here )

initialize parameters $w = 0$
for $i = 1 \cdots N$
      $c_j = \underset{t \in T}{\operatorname{argmax}} \; w^T \phi(x_i, t)$
      if $c_j \neq t_i$ then
          $w = w + \phi(x_i, t_i) - \phi(x_i, c_j)$

w is loop invariant and is the weighted average

$w = \sum_{ij} \alpha_{ij} \left( \phi(x_i, t_i) - \phi(x_i, c_j) \right)$

$f(x) = w^T \phi(x, t) = \sum_{ij} \alpha_{ij} \left( \phi(x_i, t_i)^T \phi(x, t) - \phi(x_i, c_j)^T \phi(x, t) \right)$

① KP

define $f(x) = w^T x = \sum_n \alpha_n t_n K(x_n, x)$
initialize dual parameter $\alpha_n = 0$
for $i = 1 \cdots N$
      $h_n = sign(f(x))$     ⎱ Repeat
      if $h_n \neq t_i$ then
          $\alpha_n \alpha_{n+1}$          TEST:  $y(x) = sign(f(x))$

$$MI(X,Y) = \sum_x \sum_y p(x,y) \cdot \ln \frac{p(x,y)}{p(x)\cdot p(y)} = \begin{cases} 0 & \text{when } x,y \text{ Indp} \\ \max & \text{when } x \approx y \end{cases}$$

⑥ mutual information
- → works with nomial
- → biased toward high nanty feature
- → may choose redundant feature

② chi-square

$$E_{ij} = \frac{N_{x=i} \cdot N_{y=j}}{N}$$

$$\chi_0^2 = \sum_{i,j} \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

$\xrightarrow{\text{l features}}$

$$K \text{ examples} \begin{vmatrix} 1 \cdots j \cdots l \\ i \cdots O_{ij} \cdots \} N_{x=i} \\ K \quad \vdots \\ \quad N_{y=j} \end{vmatrix}$$

③ pearson corr coeff

$$\rho(x,y) = \frac{E[(x-\mu_x)(y-\mu_y)]}{\sigma_x \sigma_y} = \frac{cov(x,y)}{\sigma_x \sigma_y}$$

④ SNR

$$\mu(x,y) = \frac{|\mu_+ - \mu_-|}{\sigma_x + \sigma_-}$$
(binary)

⑤ T-Test

$$T(x,y) = \frac{|\mu_+ - \mu_-|}{\sqrt{\frac{\sigma_+^2}{\mu_+} + \frac{\sigma_-^2}{\mu_-}}}$$

# Filter vs Wrapper

| Filter | Wrapper |
|---|---|
| 1. much faster since no need to train the model | • computationally expen. |
| 2. use statistical method of evaluation | • uses cross-validation |
| 3. might fail to find best subset | • finds best subset |
| 4. less prone to overfitting | • more prone to overfitting |

① SVM for regression

$$\min \; J(w,b) = \tfrac{1}{2}\|w\|^2 + C\sum_{n=1}^{N} (\xi_n + \hat{\xi}_n)$$

s.t. $\quad t_n \le t_n(w^T\phi(x_n)+b) + \varepsilon + \xi_n$

$\qquad t_n \ge t_n(w^T\phi(x_n)+b) - \varepsilon - \hat{\xi}_n$

$\qquad \xi_n, \hat{\xi}_n \ge 0 \quad + \; 1 \le n \le N$

② SVM for ranking

$$\min \; J(w,b) = \tfrac{1}{2}\|w\|^2 + C\sum \xi_{k,i,j}$$

s.t. $\quad w^T\phi(q_k, d_i) \ge w^T(\phi q_k, d_j) + 1 - \xi_{k,i,j}$

$\qquad \xi_{k,i,j} \ge 0$

# MDP markov decision process

① policy

$\pi^* = \underset{\pi}{\arg\max} \; E\left[ \sum_{t=0}^{\infty} \gamma^t R(s_t) \mid \pi \right]$

(policy expectation)

$\pi^A = \underset{a}{\arg\max} \sum_{s'} T(s,a,s') \, U(s')$

(best action policy)

$R$ = Reward

$\gamma^t R$ = discounted Reward

argmax

② utility

$U(s) = E\left[ \sum_{t=0}^{\infty} \gamma^t R(s_t) \mid \pi, \; s_0 = s \right]$

(expectation)

expectation & max only

$U(s) = R(s) + \gamma \max_{a} \sum_{s'} T(s,a,s') \, U(s')$ ←Bellman eqⁿ

* Naive Bias

3 Boolean input vectors $x_1 \, x_2 \, x_3$ and output $y$

$p(y=1)$

- # of parameters $= 2M+1 = 2 \times 3 + 1 = 7$

  $p(x_1 = 1 | y = 0) \quad p(x_1 = 1 | y = 1)$
  $p(x_2 = 1 | y = 0) \quad p(x_2 = 1 | y = 1)$
  $p(x_3 = 1 | y = 0) \quad p(x_3 = 1 | y = 1)$

- # of parameters if no conditional independence $= 1 + 2(2^M - 1) = 1 + 2(2^3 - 1)$

  $= 1 + 2(2^3 - 1) = 1 + 14 = 15$

# Example

| Train | DOC | words $(w_i)$ | class | Total |
|---|---|---|---|---|
| | 1 | Chinese Beiging Chinese | c1 | |
| 3 documents for c1 | 2 | Chinese Chinese shanghai | c1 | $n_1 = 8$ |
| | 3 | Chinese macao | c1 | |
| 1 document for c2 | 4 | Tokyo Japan Chinese | c2 | $n_2 = 3$ |
| | | | ? | 11 words |
| | | | | 6 unique |
| Test | 5 | | | |
| Total | $|D| = |D_1| + |D_2| = 3+1 = 4$ Chinese, Beiging, shanghai, macao, Tokyo Japan | | |

(left margin list)
1 chinese 1+1+1
2 Beiging 1
3 shanghai 1  } ignore
4 macao 1
5 Tokyo 1
6 Japan 1

① Vocabulary, $V = \{$ chinese, Beiging, shanghai, macao, Tokyo Japan $\}$

$$|V| = 6$$

② 

### category c1 = c

- prior $P(c_1) = \dfrac{|D_1|}{|D|} = \dfrac{3}{4}$

- # of words in class c1, $n_1 = 8$

conditional probabilities:

$p(w_1 | c_1) = P(\text{chinese} | c) = \dfrac{5+1}{8+6} = \dfrac{6}{14} = \dfrac{3}{7}$  

$n_{k'} = n_{11} = 5 \rightarrow$ Laplace smoothing  
$n_k \rightarrow |V|$

$P(w_2 | c_1) = P(\text{Beiging} | c) = \dfrac{1+1}{14} = \dfrac{2}{7}$

$P(w_3 | c_1) = P(\text{shanghai} | c) = 2/7 = P$

$P(w_2 | c_1) = P(\text{Tokyo} | c) = \dfrac{0+1}{8+6} = \dfrac{1}{14}$

$P(w_3 | c_1) = P(\text{Japan} | c) = \dfrac{0+1}{8+6} = \dfrac{1}{14}$

### category (c2 = j (japan))

- prior $p(c_2) = \dfrac{|D_2|}{|D|} = \dfrac{1}{4}$

- # of words in class c2, $n_2 = 3$

$P(w_1 | c_2) = P(\text{chinese} | j) = \dfrac{1+1}{3+6} = \dfrac{2}{9}$

$P(w_2 | c_2) = P(\text{Tokyo} | j) = \dfrac{1+1}{3+6} = \dfrac{2}{9}$

$P(w_3 | c_2) = P(\text{Japan} | j) = \dfrac{1+1}{3+6} = \dfrac{2}{9}$

③ choosing a class

prior  chinese  Japan  Tokyo

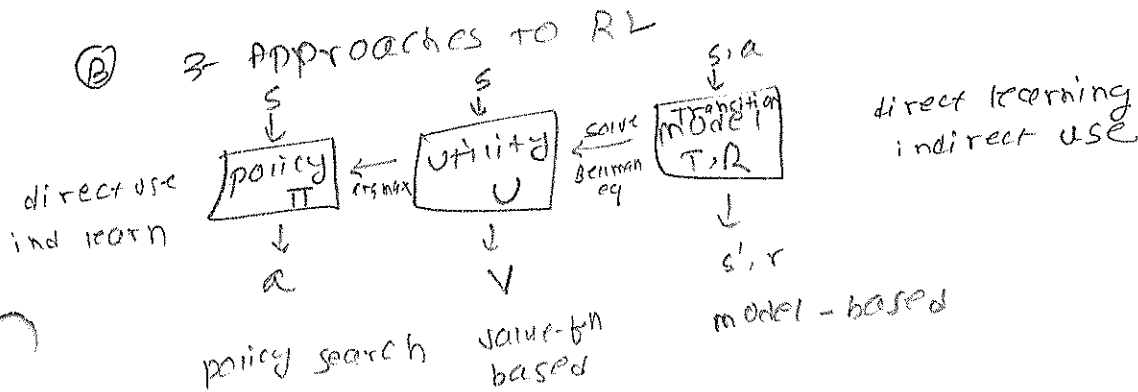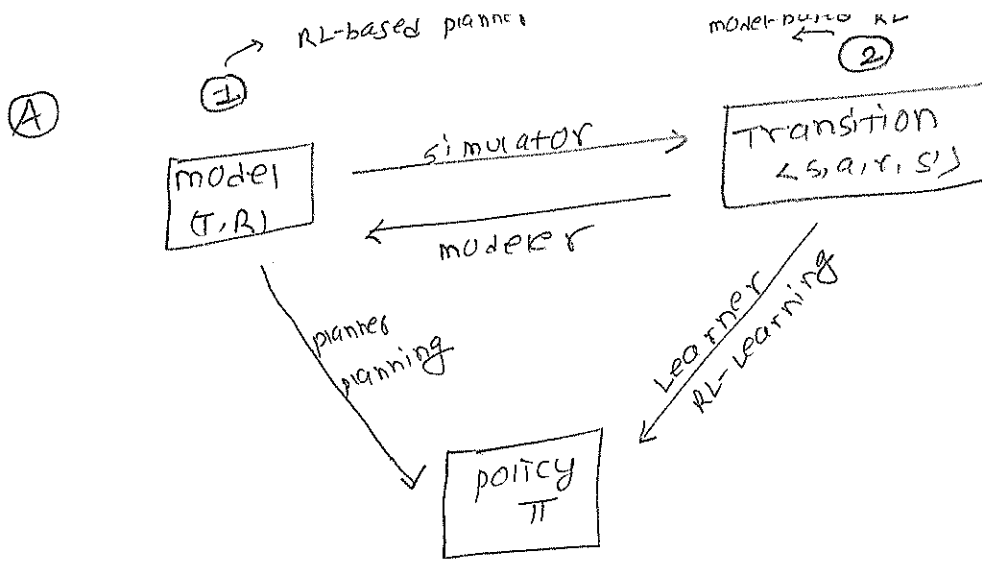$p(c_1 | x) \propto p(c_1) \prod_i p(w_i | c_1) = \dfrac{3}{4} \cdot \left(\dfrac{3}{7}\right)^3 \cdot \dfrac{1}{14} \cdot \dfrac{1}{14} \approx 0.0003$

$p(c_2 | x) \propto p(c_2) \prod_i p(w_i | c_2) = \dfrac{1}{4} \cdot \left(\dfrac{2}{9}\right)^3 \cdot \dfrac{2}{9} \cdot \dfrac{2}{9} \approx 0.0001$

$c^* = \underset{c_k}{\arg\max} \; p(c_k) \prod_{j=1}^{n} p(w_j | c_k) = \underset{c_k}{\arg\max} \{ 0.0003, 0.0001 \}$

$\rightarrow \text{largest } 0.0003 = c_1$

(A)

**model (T, R)** ⟶ *simulator* ⟶ **Transition $\langle s, a, r, s' \rangle$**

**model (T, R)** ⟵ *modeler* ⟵ **Transition**

planner / *planning* ⟶ **policy $\pi$**

*Learner / RL-learning* ⟶ **policy $\pi$**

---

(B)   3 Approaches to RL

direct learning / indirect use

direct use / and learn

$s$ ↓ **policy $\pi$** ←argmax— $s$ ↓ **utility $U$** ←solve Bellman eq— $s, a$ ↓ **Transition model T, R**

↓ $a$     ↓ $V$     ↓ $s', r$

policy search    value-f$^n$ based    model-based

---

(C)   Q function

$$U(S) = R(S) + \gamma \max_a \sum_{s'} T(s, a, s') \, U(s') \quad \text{(Bellman eqn)}$$

utility is a scalar

change U to Q

$$\pi(S) = \arg\max_a \sum_{s'} T(s, a, s') \, U(s')$$

(policy gives an action)

$$Q(s, a) = R(S) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q(s', a')$$

quiz ⟹ $$U(S) = \max_a Q(s, a)$$

⟹ $$\pi(S) = \arg\max_a Q(s, a)$$

① mutual information

$$MI(x,y) = \sum_x \sum_y p(x,y) \cdot \ln\left(\frac{p(x,y)}{p(x)\,p(y)}\right) = 0 \text{ when } x,y \text{ indp}$$

$$= KL\left[p(x,y) \,\|\, p(x)\,p(y)\right] \qquad = \text{max when } x=y$$

bad → biased towards high arity features

bad → may choose redundant feature

good → works only with nominal features + labels

# 3 parametric Approaches to

## ① Discriminant function

{ Fisher's linear dsc
  perceptron
  SVM

inference and decision
are combined as
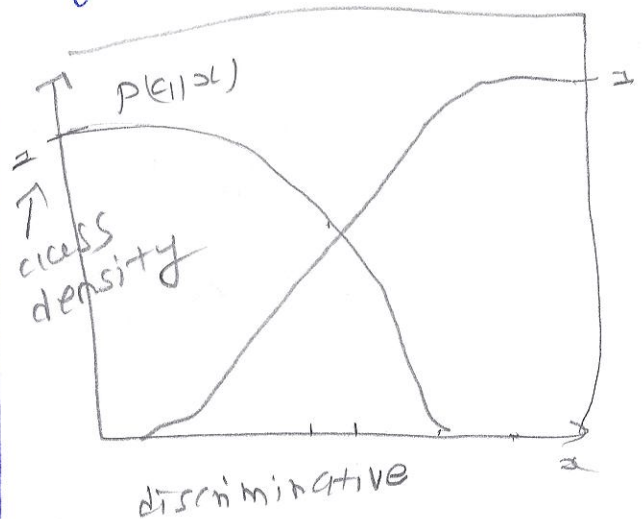single learning
problem

## ② Probabilistic ~~discriminative~~ (discriminative) models

{ logistic regression
  conditional random field

- in to dec → separate

- alters data need to
  compute $P(t_k|x)$
  than $P(t_k)(x)$

- can accommodate
  many overlapping
  features



$P(c_1|x)$

class density

discriminative

③ pearson corr coeff

$$\rho(x,y) = \frac{E[(X-\mu_x)(Y-\mu_y)]}{\sigma_x \; \sigma_y} \qquad (\text{population})$$

$$\rho(x,y) = \frac{\sum_{i=1}^{n}(x_i-\bar{x})(y_i-\bar{y})}{\sqrt{\sum_{i=1}^{n}(x_i-\bar{x})^2} \cdot \sqrt{\sum_{i=1}^{n}(y_i-\bar{y})^2}} \qquad (\text{sample})$$

Prove $\;-1 \le \rho(x,y) \le 1$

④ SNR

binary classes = $\{+,-\}$
mean = $\{\mu_+, \mu_-\}$  examples are binary
$y \in \{+1, -1\}$

$$A(x,y) = \frac{|\mu_+ - \mu_-|}{\sigma_+ + \sigma_-}$$

$\mu_+, \sigma_+ = $ mean & std
for the class
samples

⑤ T-test

$$T(x,y) = \frac{|\mu_+ - \mu_-|}{\sqrt{\dfrac{\sigma_+^2}{\mu_+} + \dfrac{\sigma_-^2}{\mu_-}}}$$

# CS 4900/5900: Machine Learning
## Fall 2017

**Class Meetings:** Tue, Thu 10:30–11:50am, ARC 212
**Instructor:** Razvan Bunescu
**Office:** Stocker 341
**Office Hours:** Tue, Thu 12:00–12:30pm, or by email appointment
**Email:** bunescu @ ohio edu
**Class Website:** http://ace.cs.ohio.edu/~razvan/courses/ml4900

**Prerequisites:**

The students are expected to be comfortable with programming and familiar with basic concepts in linear algebra and statistics.

**Textbook:**

There is no textbook for this class. Slides and supplementary materials will be made available on the course website.

**Supplementary Texts:**

> *Machine Learning: The Art and Science of Algorithms that Make Sense of Data*
> > by Peter Flach, Cambridge University Press, 2012
> *A Course in Machine Learning* [free online]
> > by Hal Daume III
> *Machine Learning*
> > by Tom Mitchell. McGraw Hill, 1997
> *Pattern Recognition and Machine Learning*
> > by Christopher Bishop. Springer, 2007
> *Pattern Classification*
> > by Richard O. Duda, Peter E. Hart, & David G. Stork. Wiley-IS, 2001
> *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*
> > by T. Hastie, R. Tibshirani, & J. H. Friedman. Springer Verlag, 2009

**Course Description:**

This course will give an overview of the main concepts, techniques, and algorithms underlying the theory and practice of machine learning. The course will cover the fundamental topics of classification, regression and clustering, and a number of corresponding learning models such as perceptrons, logistic regression, linear regression, Naive Bayes, nearest neighbors, and Support Vector Machines. The description of the formal properties of the algorithms will be supplemented with motivating applications in a wide range of areas including natural language processing, computer vision, bioinformatics, and music analysis. The topics covered in this course will also prepare students for taking more advanced courses in data mining and deep learning.

## Grading:
50%: Homework Assignments
20%: Midterm Exam (Oct 12, in class) 1 h r 20 min
30%: Final Exam (Dec 12, 10:10am – 12:10pm)

## Grading Scale:
A (> 92%) A–(> 90%) B+(> 87%) B(> 83%) B–(> 80%)
C+(> 77%) C(> 73%) C–(> 70%) D+(> 67%) D(> 63%) D–(> 60%)

## Important Dates:
Friday, Sep 1: Last day to add class.
Tuesday, Oct 10: Reading Day, no class.
Friday, Nov 3: Last day to drop class.
Thursday, Nov 23: Thanksgiving break, no class.
Thursday, Dec 7: Last day of this class.

## Course and Attendance policies:
**Assignments:** All homework assignments are due before the class. No late submissions will be accepted without prior approval.

**Attendance:** It is in your best interest to attend the lectures. Some of the material will not be found in the supplementary text or on the slides. Extra credit will be awarded for class activity. Also, be sure to check your OU email for important announcements on a regular basis.

## Academic Dishonesty Policy:
All work must be the student's own. All external references used in reports must be properly cited. No credit will be given for duplicate or plagiarized work. Additional measures may be imposed by the University Judiciaries, when conditions warrant. Students may appeal academic sanctions through the grade appeal process. The OU Student Code of Conduct Policy is available online at:

http://www.ohio.edu/communitystandards/academic/students.cfm

## Disability-based Accommodation:
Any student who suspects s/he may need an accommodation based on the impact of a disability should contact the class instructor privately to discuss the students specific needs and provide written documentation from the Office of Student Accessibility Services. If the student is not yet registered as a student with a disability, s/he should contact the Office of Student Accessibility Services.

## Other Policies:
Be sure to notify the professor of any exam conflicts or other extenuating circumstances well in advance. No missed exams will be made up without prior approval. Medical excuse forms need to explicitly mention that the student could not have attended the exam at the specified time due to health concerns.