

140

HW Assignment 1 (Due by 10:30am on Sep 26)

Bhishan Poudel

1 Theory (40 points)

1. [Polynomial Curve Fitting, 20 points]

Consider the problem of fitting a dataset of N points with a polynomial of degree M , by minimizing the sum-of-squares error:

$$J(\mathbf{w}) = \frac{1}{2N} \sum_{n=1}^N (h_{\mathbf{w}}(\mathbf{x}_n) - t_n)^2 \quad (1)$$

where $h_{\mathbf{w}}(\mathbf{x}) = \sum_{j=0}^M w_j x^j$. We have shown in class that the solution to minimizing $J(\mathbf{w})$ satisfies the following set of linear equations:

$$\sum_{j=0}^M A_{ij} w_j = T_i \quad (2)$$

$$\text{where } A_{ij} = \sum_{n=1}^N x_n^{i+j} \text{ and } T_i = \sum_{n=1}^N x_n^i t_n \quad (3)$$

Derive the solution for the regularized version of polynomial curve fitting, which minimizes the objective function below:

$$J(\mathbf{w}) = \frac{1}{2N} \sum_{n=1}^N (h_{\mathbf{w}}(\mathbf{x}_n) - t_n)^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2 \quad (4)$$

2. [Probability Theory, 20 points] *Exercise 1.3, page 58 in PRML.*

Suppose that we have three coloured boxes r (red), b (blue), and g (green). Box r contains 3 apples, 4 oranges, and 3 limes, box b contains 1 apple, 1 orange, and 0 limes, and box g contains 3 apples, 3 oranges, and 4 limes. If a box is chosen at random with probabilities $p(r) = 0.2$, $p(b) = 0.2$, $p(g) = 0.6$, and a piece of fruit is removed from the box (with equal probability of selecting any of the items in the box), then what is the probability of selecting an apple? If we observe that the selected fruit is in fact an orange, what is the probability that it came from the green box?

3. [Ridge Regression (*), 20 points]

Consider the regularized linear regression objective shown below:

$$J(\mathbf{w}) = \frac{1}{2N} \sum_{n=1}^N (h(\mathbf{x}_n, \mathbf{w}) - t_n)^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

- (a) Minimizing the L2 norm of \mathbf{w} drives all parameters, including w_0 , towards 0. Are there situations in which we do not want to constrain w_0 to be small? If yes, give an example, if not show why it is useful to constrain all the weights to be small, including w_0 .
- (b) We have seen in class how to compute the weights \mathbf{w} that minimize $J(\mathbf{w})$. Assume now that we replace $\|\mathbf{w}\|$ in $J(\mathbf{w})$ with $\|\mathbf{w}_{[1:]}\|$, where $\mathbf{w}_{[1:]} = [w_1, w_2, \dots, w_M]$. Derive the solution for \mathbf{w} that minimizes this new objective function.

2 Implementation (80 points)

In this exercise, you are asked to run an experimental evaluation of linear regression on the Athens houses dataset, and on an artificial dataset with and without L2 regularization. The input data is available at <http://ace.cs.ohio.edu/~razvan/courses/ml4900/hw01.zip>. Make sure that you organize your code in folders as shown in the table below. Write code only in the Python files indicated in bold.

$w = (X^T X)^{-1} e^T e + t = \frac{\text{two } w_1, w_2, w_3, t}{\text{row 22 array}}$ $\text{RMSE} = \sqrt{\frac{\sum_{n=1}^N (h_n - t_n)^2}{N}}$ $= \frac{\text{np.sqrt(np.sum((h-t)^2))}}{N}$ $= \text{np.sqrt}((h-t)^2 \cdot \text{mean}())$	<pre>ml4900/ hw01/ code/ univariate.py multivariate.py polyfit.py train_test_line.png data/ univariate/ train.txt, test.txt multivariate/ train.txt, test.txt polyfit/ train.txt, test.txt, devel.txt</pre>
---	---

1. [Univariate Regression, 20 points]

Train a univariate linear regression model to predict house prices as a function of their floor size, based on the solution to the system with 2 linear equations discussed in class. Use the dataset from the folder `hw01/data/univariate`. Python3 skeleton code is provided in `univariate.py`. After training print the parameters and report the RMSE and the objective function values on the training and test data. Plot the training using the default blue circles and test examples using lime green triangles. On the same graph also plot the linear approximation.

2. [Multivariate Regression, 20 points]

Train a univariate linear regression model to predict house prices as a function of their floor size, number of bedrooms, and year. Use the normal equations discussed in class, and evaluate on the dataset from the folder `hw01/data/multivariate`. Python3 skeleton code is provided in `multivariate.py`. After training print the parameters and report the RMSE and the objective function values on the training and test data. Compare the test RMSE with the one from the univariate case above.

3. [Polynomial Curve Fitting, 40 points]

In this exercise, you are asked to run an experimental evaluation of a linear regression model, with and without regularization. Use the normal equations discussed in class, and evaluate on the dataset from the folder `hw01/data/polyfit`.

- (a) Select 30 values for $x \in [0, 1]$ uniformly spaced, and generate corresponding t values according to $t(x) = \sin(2\pi x) + x(x+1)/4 + \epsilon$, where $\epsilon = N(0, 0.005)$ is a

zero mean Gaussian with variance 0.005. Save and plot all the values. Done in `dataset.txt`.

- (b) Split the 30 samples (x_n, t_n) in three sets: 10 samples for training, 10 samples for validation, and 10 samples for testing. Save and plot the 3 datasets separately. Done in `train.txt`, `test.txt`, `devel.txt`.
- (c) Consider a linear regression model with polynomial basis functions, trained with the objective shown below:

$$J(\mathbf{w}) = \underbrace{\frac{1}{2N} \sum_{n=1}^N (h(x_n, \mathbf{w}) - t_n)^2}_{\text{hat MSE}} + \underbrace{\frac{\lambda}{2} \|\mathbf{w}\|^2}_{\text{L2 Regularization}}$$

$\sum_{j=1}^M w_j$
generally we don't
regularize w_0 .

Show the closed form solution (vectorized) for the weights \mathbf{w} that minimize $J(\mathbf{w})$.

- (d) Train and evaluate the linear regression model in the following scenarios:
1. Without regularization: Use the training data to infer the parameters \mathbf{w} for all values of $M \in [0, 9]$. For each order M , compute the RMSE separately for the training and test data, and plot all the values on the same graph, as shown in class.
 2. With regularization: Fixing $M = 9$, use the training data to infer the parameters \mathbf{w} , one parameter vector for each value of $\ln \lambda \in [-50, 0]$ in steps of 5. For each parameter vector (lambda value), compute the RMSE separately for the training and validation data, and plot all the values on the same graph, as shown in class. Select the regularization parameter that leads to the parameter vector that obtains the lowest RMSE on the validation data, and use it to evaluate the model on the test data. Report and compare the test RMSE with the one obtained without regularization.

3 Submission

Turn in a hard copy of your homework report at the beginning of class on the due date. Electronically submit on Blackboard a `hw01.zip` file that contains the `hw01` folder in which your code is in the 3 required files.

On a Linux system, creating the archive can be done using the command:

```
> zip -r hw01.zip hw01.
```

Please observe the following when handing in homework:

1. Structure, indent, and format your code well.
2. Use adequate comments, both block and in-line to document your code.
3. On the theory assignment, **clear and complete explanations and proofs of your results are as important as getting the right answer.**
4. Make sure your code runs correctly when used in the directory structure shown above.

HW 1

Due: Sep 19 TU

Bhishan Poudel

Q1.1 polynomial curve fitting

cost function without the regularization

is

$$J(w) = \frac{1}{2N} \sum_{n=1}^N (h_w(x_n) - t_n)^2 \quad \textcircled{1}$$

where predicted values

$$h_w(x) = \underbrace{w_0 x_0}_{=w_0} + w_1 x_1 + w_2 x_2 + \dots + w_M x_M$$

w_0, w_1, \dots, w_M
= model parameters
 w_0 = bias term or
intercept term

$$h_w(x) = \sum_{j=0}^M w_j x_j$$

minimizing the cost function w/o regularization,

$$\frac{\partial J}{\partial w} = 0 \quad \text{gives}$$

$$\sum_{j=0}^M A_{ij} w_j = T_i \quad \textcircled{2}$$

where

$$A_{ij} = \sum_{n=1}^N x_n^{i+j}$$

$$T_i = \sum_{n=1}^N x_n^i + t_n$$

\textcircled{3}

problem: minimize the function, λ is shrinkage

$$J(\mathbf{w}) = \frac{1}{2N} \sum_{n=1}^N \underbrace{\left(h_{\mathbf{w}}(\mathbf{x}_n) - t_n \right)^2}_{\frac{1}{2} \cdot \text{SSE}} + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

L2 regularizer
or
shrinkage penalty

solution: the hypothesis is

$$h_{\mathbf{w}}(\mathbf{x}_n) = \sum_{j=0}^M w_j x_n^j$$

$n = 0, \dots, N \rightarrow N \text{ samples } x_1, x_2, \dots, x_N$
 $j = 0, \dots, M \rightarrow M \text{ features } w_0, w_1, \dots, w_M$
 bias 特征 age
 w_0 w_1 w_2
 $t = \text{target (this)} \quad (\text{price})$

we can write cost function as,

$$J(\mathbf{w}) = \frac{1}{2N} \sum_{n=1}^N \left(\sum_{j=0}^M w_j x_n^j - t_n \right)^2 + \frac{\lambda}{2} \sum_{j=0}^M w_j^2$$

we want to minimize the cost function $J(\mathbf{w})$ w.r.t. weights w_i ,

$$\nabla_{w_i} J(\mathbf{w}) = 0$$

$$\frac{\partial J(w_i)}{\partial w_i} = 0$$

$$0 = \frac{1}{2N} \cdot 2 \sum_{n=1}^N \left(\sum_{j=0}^M w_j x_n^j - t_n \right) \cdot x_n^i + \frac{\lambda}{2} \cdot 2 w_i$$

$$0 = \sum_{n=1}^N \sum_{j=0}^M w_j x_n^{ij} - \underbrace{\sum_{n=1}^N t_n x_n^i}_{\text{bias}} + \lambda N w_i$$

$$\text{Or, } \sum_{j=0}^M \left(\sum_{n=1}^N \alpha_n^{j+n} \right) w_i + \gamma_N w_i = \sum_{n=1}^N t_n \omega_n$$

$$\Rightarrow \boxed{\sum_{j=0}^M A_{ij} w_j + \gamma_N w_i = T_i} \quad \# \text{ Answer}$$

where,

$$A_{ij} = \sum_{n=1}^N \alpha_n^{i+j}$$

$$T_i = \sum_{n=1}^N t_n \omega_n^i$$

(2)
✓

QNT2 probability Theory By Bhopal

Ex 1.3, page 58 in PRML

Red R: 3a 4o 3l

blue B: 2a 2o 0l

green G: 3a 3o 4l

$$P(R) = 0.2 \quad P(B) = 0.2 \quad P(G) = 0.6$$

- If a piece of fruit is removed from the box, then what is the prob of selecting an apple?

- If we observe that the selected fruit is an apple, what is the prob that it came from the green box?

Solution:

⑥	3a 4o 3l
---	----------------

Red

$$N_R = 10$$

2a 2o 0l

Blue

$$N_B = 2$$

3a 3o 4l

Green

$$N_G = 10$$

$$\begin{aligned} N &= N_R + N_B + N_G \\ &= 10 + 2 + 10 \\ &= 22 \end{aligned}$$

$$P(A|R) = \frac{3}{10} = 0.3 \quad P(A|B) = \frac{1}{2} = 0.5 \quad P(A|G) = \frac{3}{10} = 0.3 \quad \left. \right\} \text{from figure}$$

$$P(O|R) = \frac{4}{10} = 0.4 \quad P(O|B) = \frac{1}{2} = 0.5 \quad P(O|G) = \frac{3}{10} = 0.3$$

$$P(R) = 0.2 \quad P(B) = 0.2 \quad P(G) = 0.6 \quad] \text{given}$$

Implementation

QN 3.C

Given objective function for linear regression with polynomial basis functions is,

(3)

Cost function

$$J(\vec{w}) = \frac{1}{2N} \sum_{n=1}^N [h(x_n; \vec{w}) - t_n]^2 + \frac{\lambda}{2} \|\vec{w}\|^2$$

Now, we have to find the vectorized solution of w which minimizes cost function J .

For this,

we first write hypothesis in terms of design matrix X and weight vector w .

$$h = Xw \quad \text{--- (1)}$$

Then, eqn (1) can be written as,

$$J = \frac{1}{2N} (Xw - t)^T (Xw - t) + \frac{\lambda}{2} w^T w$$

We minimize this equation w.r.t. w ,

$$\frac{\partial J}{\partial w}$$

$$= \frac{1}{2N} \cdot 2X^T(Xw - t) + \frac{\lambda}{2} \cdot 2w$$

more we have used matrix derivative

for formula

$$\frac{d}{d\omega} (x^T x) = 2x$$

$$\frac{d}{d\omega} (A^T \omega + b) = 2 A^T (A \omega + b)$$

Ansible matrices

formula

then,

$$0 = x^T x w - x^T t + d N M$$

$$x^T t = (x^T x + d N) w$$

$$w = (d N I + x^T x)^{-1} (x^T t) \quad ||$$

ANSWER

Note that $(x^T x)^{-1} x^T$ is called Moore-Penrose
pseudoinverse of any shape matrix X .

e.g. $x^T x = \begin{pmatrix} 4 & 5 & 0 \\ 5 & 8 & 0 \\ 0 & 0 & 4 \end{pmatrix} = \begin{pmatrix} 4 & 4 \\ 4 & 4 \end{pmatrix} = I_{4,4}$

$$x^T t = \begin{pmatrix} 4 & 5 & 0 \\ 5 & 8 & 0 \\ 0 & 0 & 1 \end{pmatrix} = \begin{pmatrix} 4 & 1 \\ 4 & 1 \end{pmatrix}$$

$$w = \begin{pmatrix} 4 & 1 \\ 4 & 1 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

$$w^T = [w_0 \ w_1 \ w_2 \ w_3] = x^{(n+1)}$$

$$I_{4,4} = \begin{bmatrix} 1 & & 0 & \\ & 1 & & 0 \\ 0 & & 1 & \\ & & & 1 \end{bmatrix}$$

Set $I^{(0)}[0] = 0$
 ~~$I^{(0)}[1] = 1$~~

we don't regularize bias term.

$$= \begin{bmatrix} 0 & & 0 & \\ & 1 & & 0 \\ 0 & & 1 & \\ & & & 1 \end{bmatrix}$$

if we regularize bias term, it will shift the mean of target vector.

Then adding const value to all the target DOES NOT hypothesis similarity shifted hypothesis.

I-1 (20)

I-2 (20)

I-3 (40)

HW Assignment 3 (Due by 10:30am on Oct 12)

(106/100)

(120)

1 Theory (100 points)**1. [Maximum Likelihood, 20 points]**

The Poisson distribution specifies the probability of observing k events in an interval, as follows:

$$P(k \text{ events in interval}) = e^{-\lambda} \frac{\lambda^k}{k!} \quad (1)$$

For example, k can be the number of meteors greater than 1 meter diameter that strike Earth in a year, or the number of patients arriving in an emergency room between 10 and 11 pm¹.

Suppose we observe N samples k_1, k_2, \dots, k_N from this distribution (i.e. numbers of meteors that strike Earth over a period of N years). Derive the maximum likelihood estimate of the event rate λ .

2. [Logistic Regression, 20 points]

Consider a dataset that contains the 4 examples below i.e., the truth table of the logical XOR function. Prove that no logistic regression model can perfectly classify this dataset. Do not forget the bias feature $x_0 = 1$.

x_1	x_2	t
0	0	0
0	1	1
1	0	1
1	1	0

Hint: Prove that there cannot be a vector of parameters \mathbf{w} such that $P(t = 1|\mathbf{x}, \mathbf{w}) \geq 0.5$ for all examples \mathbf{x} that are positive, and $P(t = 1|\mathbf{x}, \mathbf{w}) < 0.5$ for all examples \mathbf{x} that are negative.

3. [Logistic Regression, 20 points]

Prove that the gradient (with respect to \mathbf{w}) of the negative log-likelihood error function for logistic regression corresponds to the formula shown in lecture 4:

$$\nabla_{\mathbf{w}} E(\mathbf{w}) = \sum_{n=1}^N (h_n - t_n) \mathbf{x}_n \quad (2)$$

4. [Logistic Regression, 20 points]

In `scikit`, the objective function for logistic regression expresses the trade-off between training error and model complexity through a parameter C that is multiplied with the error term, as shown below. See the `scikit` documentation at http://scikit-learn.org/stable/modules/linear_model.html#logistic-regression.

$$E(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C * \sum_{n=1}^N \ln(e^{-t_n(\mathbf{w}^T \mathbf{x}_n)} + 1) \quad (3)$$

¹https://en.wikipedia.org/wiki/Poisson_distribution

- Show that the sum in the second term is equal with the negative log-likelihood, where $t_n = +1$ stands for positive labels and $t_n = -1$ stands for negative labels.
 - Compute the C parameter such that the objective is equivalent with the standard formulation shown on the slides in which the regularization parameter λ is multiplied with the L2 norm term.
5. [Softmax Regression, 20 points]
 Show that Logistic Regression is a special case of Softmax Regression. That is to say, if w_1 and w_2 are the parameter vectors of a Softmax Regression model for the case of two classes, then there exists a parameter vector w for Logistic Regression that results in the same classification as the Softmax Regression model.
6. [Softmax Regression (*), 20 points]
 Prove that the gradient (with respect to w_k) of the negative log-likelihood error function for regularized softmax regression corresponds to the formula shown in lecture 4, for any class $k \in [1..K]$:

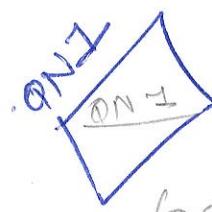
$$\nabla_{w_k} E(w) = -\frac{1}{N} \sum_{n=1}^N (\delta_k(t_n) - p(C_k | \mathbf{x}_n)) \mathbf{x}_n + \alpha w_k \quad (4)$$

2 Submission

Turn in a hard copy of your homework report at the beginning of class on the due date. On this theory assignment, **clear and complete explanations and proofs of your results are as important as getting the right answer.**

HW 03

Bhishan Poudel

maximum LikelihoodDerive MLE of event rate parameter λ)

Let x_1, x_2, \dots, x_n be iid poisson random variables with prob mass function,

pmf of poisson distn

$$p(x_i | \lambda) = \frac{e^{-\lambda} \lambda^{x_i}}{x_i!}$$

$$= p(x_i)$$

Now, the joint prob of all variables x_i , called likelihood function, is given by

likelihood

$$L(\lambda) = \prod_{i=1}^n p(x_i | \lambda)$$

$$= \prod_{i=1}^n \frac{e^{-\lambda} \lambda^{x_i}}{x_i!}$$

The log likelihood of pmf is,

Log likelihood

$$\ln L(\lambda) = \ln \prod_{i=1}^n \frac{e^{-\lambda} \lambda^{x_i}}{x_i!}$$

$$= \sum_{i=1}^n \ln \frac{e^{-\lambda} \lambda^{x_i}}{x_i!}$$

$$= \sum_{i=1}^n [-\lambda + x_i \ln \lambda - \ln x_i!]$$

$$\ln L(\lambda) = -n\lambda + n\ln \lambda - \sum_{i=1}^n \ln x_i!$$

①

To get the maximum likelihood estimate of the parameter λ , we maximize the Log Likelihood function $L(\lambda)$ w.r.t λ . [maximizing $L(\lambda)$ is same as minimizing $-n\ln L(\lambda)$]

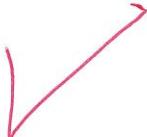
$$0 = \frac{\partial}{\partial \lambda} \ln L(\lambda)$$

$$= \frac{\partial}{\partial \lambda} [-n\lambda + n \ln \bar{x}_i - \bar{x}_i \ln \bar{x}_i]$$

$$0 = -n + \frac{\bar{x}_i}{\lambda} - 0$$

$$\lambda = \frac{\bar{x}_i}{n}$$

(25)

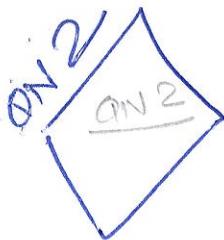


MLE of λ

$$\boxed{\lambda = \frac{\sum_{i=1}^n x_i}{n}} \rightarrow \text{number of samples}$$

Ans

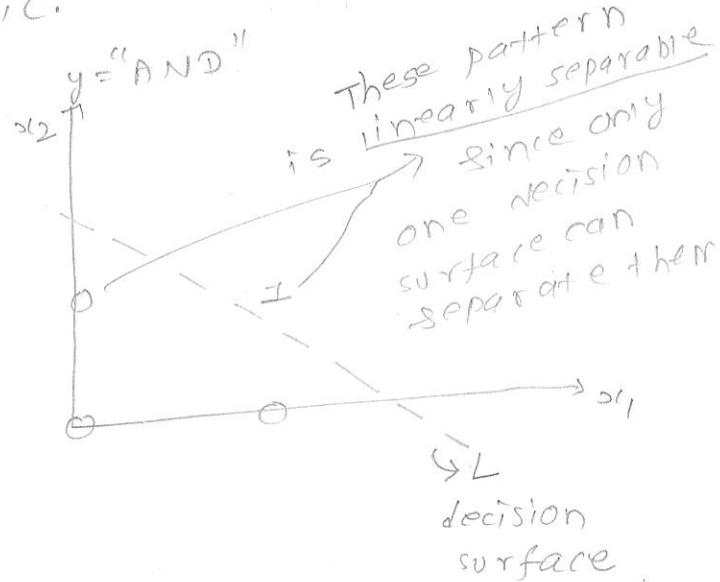
Hence, the maximum likelihood estimate of the Poisson distribution parameter λ is just the mean (or expectation) of the distribution.



XOR problem in Logistic Regression

To describe XOR problem in LR, I shall start with "AND" logic.

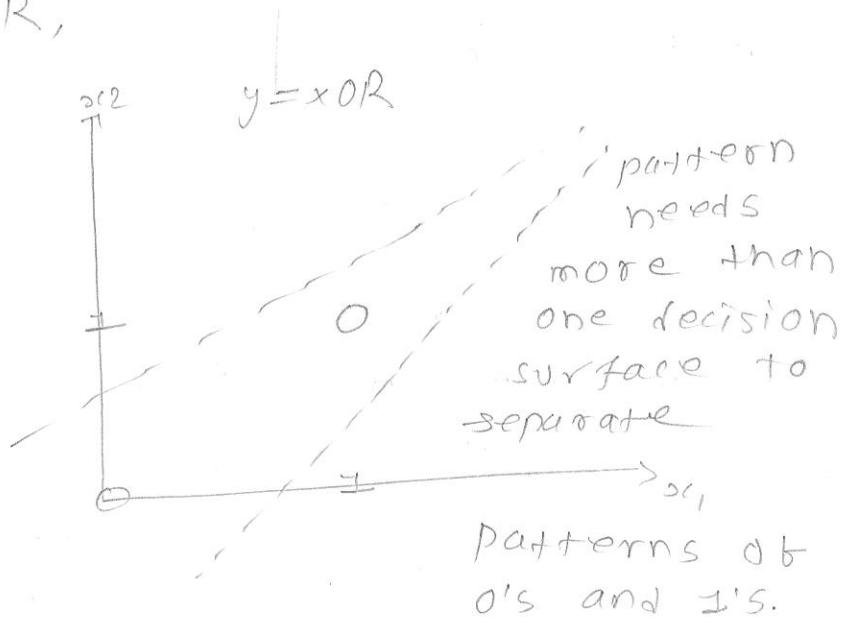
x_1	x_2	$y = x_1 \wedge x_2$
0	0	0
0	1	0
1	0	0
1	1	1



Now, look at XOR,

Truth table

x_1	x_2	$y = \text{XOR}$
0	0	0
0	1	1
1	0	1
1	1	0



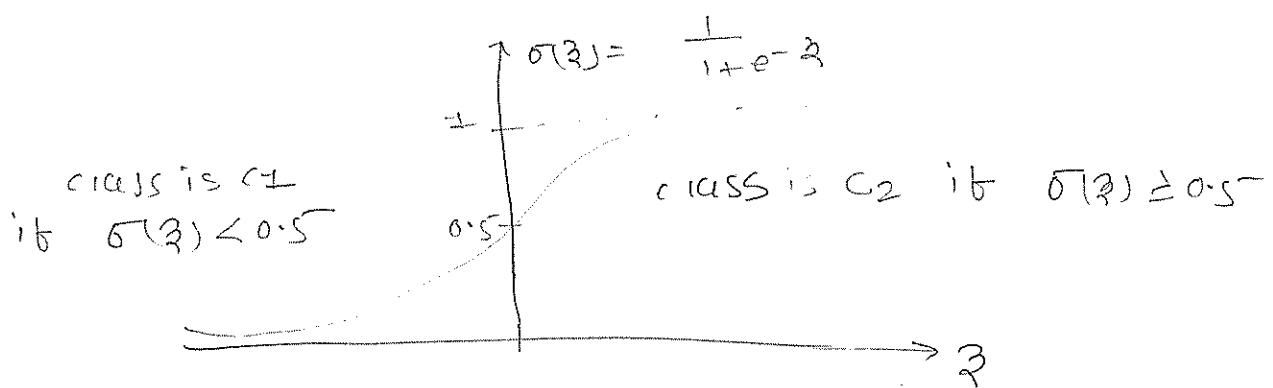
\therefore XOR logic is not linearly separable.

P.T.O.

Now, we shall show that Logistic regression is linear classifier:

The classifier used in LR is sigmoid

$$\text{function } \sigma(z) = \frac{1}{1+e^{-z}}$$



In Logistic Regression,

$$\text{pattern is } + \text{ if } \frac{1}{1+e^{-w^T x}} \geq \frac{1}{2}$$

$$0 \text{ if } \frac{1}{1+e^{-w^T x}} < \frac{1}{2}$$

let's set the classifier to separator value.

$$\frac{1}{1+e^{-w^T x}} = \frac{1}{2}$$

$$z = 1 + e^{-w^T x}$$

$$\frac{1}{z} = e^{-w^T x} = \frac{1}{e^{w^T x}}$$

$$e^{w^T x} = 1$$

$$w^T x = \ln 1 = 0$$

$$\Rightarrow \boxed{\sum_i w_i x_i = 0}$$

this means L is linear classifier.

conclusion: \Rightarrow 1) this linear classifier

2) XOR problem is non-linear and
linearly inseparable

use:

$w^T x_n > 0$ if $x_n \in +$

$w^T x_n < 0$ if $x_n \in -$

then put 4 samples into
the inequalities.

perfectly classify
dataset that
follows XOR
logic.

ANS

Better method

x_0	x_1	x_2	t
1	0	0	0
1	0	1	1
1	1	0	1
1	1	1	0

for linear separability, we want

$$w_0 + \sum_i w_i x_i < 0 \text{ if } t_i = 0$$

$$w_0 + \sum_i w_i x_i \geq 0 \text{ if } t_i = 1$$

Note,

$$w_0 < 0 \quad \text{--- (1)}$$

$$w_0 + w_2 \geq 0 \quad \text{--- (2)} \quad \left. \begin{array}{l} \\ \end{array} \right\} w_0 + (w_0 + w_1 + w_2) \geq 0$$

$$w_0 + w_1 \geq 0 \quad \text{--- (3)} \quad \text{but } w_0 < 0$$

$$w_0 + w_1 + w_2 < 0 \quad \text{--- (4)} \quad \text{and } w_0 + w_1 + w_2 < 0$$

\therefore it contradicts our assumptions of linear separability.

ON3
ON3

Gradient of Logistic Regression cost

(Bind & E)

for linear regression, hypothesis $h = w^T x$

for logistic regression, hypothesis $h = \sigma = \frac{1}{1+e^{-w^T x}}$

Aside:
binomial dist'n

$$p(x; p) = n! \cdot p^x \cdot (1-p)^{n-x}$$

for $x \in \{0, \dots, n\}$

$$n! = \frac{n!}{x!(n-x)!}$$

$$h = \sigma = \frac{1}{1+e^{-w^T x}}$$

CI write $w^T x$
as wx since they
are just matrix
dot product of two

matrices)

Binomial
Likelihood function (Bernoulli)

$$\left. \begin{array}{l} \text{prob for logistic regression} \\ p(t|w) = h^n \cdot (1-h)^{1-n} \end{array} \right\}$$

Likelihood function

$$L(w) = \prod_{n=1}^N h^n \cdot (1-h)^{1-n}$$

(Bishop eq 4.89
p. 206/223)

-ve log likelihood, E or $J = -\ln L(w)$

$$E = -\ln \prod_{n=1}^N h^n \cdot (1-h)^{1-n}$$

$$E = \sum_n E_n = \sum_{n=1}^N (-t_n \ln h_n - (1-t_n) \ln (1-h_n))$$

$$\text{Loss, } E(w) = -\frac{1}{N} \sum_n [t_n \ln h_n + (1-t_n) \ln (1-h_n)]$$

for any n th sample the cost can be written as,

LOSS

$$E = -t_n \ln h - (1-t_n) \ln (1-h) \quad \text{--- (2)}$$

-ve log
likelihood
error fn
for Vf

$$\text{here } h = \sigma = \sigma(w^T x) = \sigma(wx)$$

Before proceeding further, I would derive derivative of sigmoid function.

$$\sigma(z) = \frac{1}{1+e^{-z}} \Rightarrow \frac{d\sigma}{dz} = \frac{-1 \cdot e^{-z} \cdot (-1)}{(1+e^{-z})^2}$$

$$\begin{aligned}\frac{d\sigma}{dz} &= \frac{e^{-z}}{(1+e^{-z})^2} = \frac{1}{1+e^{-z}} \cdot \frac{(e^{-z}+1)-1}{1+e^{-z}} \\ &= \frac{1}{1+e^{-z}} \cdot \left(\frac{e^{-z}+1}{1+e^{-z}} - \frac{1}{1+e^{-z}} \right) \\ &= \frac{1}{1+e^{-z}} \cdot \left(1 - \frac{1}{1+e^{-z}} \right)\end{aligned}$$

$$\boxed{\frac{d\sigma(z)}{dz} = \sigma(1-\sigma)}$$

similarly $\boxed{\frac{d\sigma(a_2)}{dz} = \alpha \sigma(1-\sigma)}$ (here $a \neq a(z)$)

Derivative
Sigmoid
function

$$\boxed{\frac{d\sigma(a_2)}{dz} = \alpha \sigma(1-\sigma)}$$

Now, going back to the problem, let's calculate the gradient of loss function.

$$\frac{\partial E}{\partial w} = \frac{2}{n} \left[-t \ln h - (1-t) \ln(1-h) \right]$$

$$= -\frac{1}{n} \frac{\partial h}{\partial w} - \frac{(1-t)}{1-h} \cdot \frac{\partial h}{\partial w}$$

$$= -\frac{1}{n} \frac{\partial \sigma(wx)}{\partial w} + \frac{1-t}{1-h} \frac{\partial \sigma(wx)}{\partial w}$$

$$= -\frac{1}{n} \times t(1-h) + \frac{1-t}{1-h} \times h(1-h)$$

$$\therefore \frac{d\sigma(\theta)}{d\theta} = a\sigma(1-\sigma)$$

(W)

✓

$$= -t \times \cancel{-txh} + xh - \cancel{txh}$$

$$= xh - tx$$

$$\boxed{\frac{\partial E}{\partial w} = (h-t)x}$$

(here this is for a single sample, but for total loss we add up all the losses.)

$$\Rightarrow \boxed{\nabla_w E = \sum_{n=1}^N (h_n - t_n) x_n}$$

Q.E.D.

Q34
ONLY

Logistic Regression in SKLEARN

~~parts~~ the cost function for L2 penalized logistic regression is given by,

L2 Regularized cost fn
for LR
(in SKLearn)

$$E = \frac{1}{2} w^T w + C \sum_{n=1}^N \ln(1+e^{-t_n w^T x_n})$$

L2 regularizer cost for LR

where $t_n \in \{-1, 1\}$

①

we have to show that the summation

permits eval to the -ve log likelihood,

-ve log
likelihood

$$ED = -\ln p = \sum_{n=1}^N \ln(1+e^{-t_n w^T x_n}) \quad \text{--- } ②$$

Now, the posterior probabilities for class 0 and 1 are,

posterior
prob
for
LR

$$p(c_1|x) = \sigma(w^T x) = \frac{1}{1+e^{w^T x}}$$

$$p(c_0|x) = 1 - \sigma(w^T x) = \frac{e^{-w^T x}}{1+e^{-w^T x}}$$

③

The likelihood function is,

Likelihood
for
LR

$$P(t|w) = \prod_{n=1}^N h_n^{t_n} (1-h_n)^{(1-t_n)} \quad \text{--- } ④$$

11

The -ve log likelihood is,

$$\begin{aligned}-\text{un } p(\mathbf{x}|\mathbf{w}) &= -\text{un } \prod_{n=1}^N h_n^{t_n} \cdot (1-h_n)^{1-t_n} = -\sum_{n=1}^N \text{un}(h_n^{t_n}, (1-h_n)^{1-t_n}) \\&= -\sum_{n=1}^N [\text{un}(h_n^{t_n}) + \text{un}(1-h_n)^{1-t_n}]\end{aligned}$$

cost for LR

$E_D = -\text{un } p = -\sum_{n=1}^N [\underbrace{\text{un } h_n}_{\text{for class 0}} + \underbrace{(1-h_n)}_{\text{for class 1}} \text{un}(1-h_n)] \quad \textcircled{5}$

where $t_n \in \{0, 1\}$

$$h_n = \frac{1}{1+e^{-\mathbf{w}^T \mathbf{x}_n}}$$

$$1-h_n = 1 - \frac{1}{1+e^{-\mathbf{w}^T \mathbf{x}_n}}$$

$$= \frac{1+e^{-\mathbf{w}^T \mathbf{x}_n} - 1}{1+e^{-\mathbf{w}^T \mathbf{x}_n}}$$

$$= \frac{e^{-\mathbf{w}^T \mathbf{x}_n}}{1+e^{-\mathbf{w}^T \mathbf{x}_n}}$$

$$= \frac{1}{e^{\mathbf{w}^T \mathbf{x}_n} + 1}$$

$1-h_n = \frac{1}{1+e^{\mathbf{w}^T \mathbf{x}_n}}$

$$E_D = -\sum_{n=1}^N t_n p$$

$$= -\sum_{n=1}^N \left[t_n \ln \left(\frac{1}{1+e^{-w^T x_n}} \right) + (1-t_n) \ln \left(\frac{1}{1+e^{w^T x_n}} \right) \right]$$

$t_n \in \{0, 1\}$

$$= -\sum_{n=1}^N \begin{cases} \ln \frac{1}{1+e^{-w^T x_n}} & \text{if } t_n = 1 \\ \ln \frac{1}{1+e^{w^T x_n}} & \text{if } t_n = 0 \end{cases}$$

$$= -\sum_{n=1}^N t_n \ln \frac{1}{1+e^{-t_n w^T x_n}} \quad \text{if } t_n \in \{-1, 1\}$$

(F4)

$$\boxed{E_D = -\sum_{n=1}^N t_n \ln (1+e^{-t_n w^T x_n})} \quad Q.E.D.$$

$t_n \in \{-1, 1\}$



QN 4b Find 'C' parameter of LR cost fn in sklearn.

Note, the L2 regularized cost function for LR

$$E = E_w + E_D$$

$$= \frac{1}{2} w^T w + (-\ln p)$$

when the t_n 's are {0, 1}

$$E = \frac{1}{2} w^T w - \sum_{n=1}^N [t_n \ln h_n + (1-t_n) \ln (1-h_n)] \quad \text{⑥}$$

$t_n \in \{0, 1\}$

when the t_n 's are {-1, 1}

$$E = \frac{1}{2} w^T w + \sum_{n=1}^N \ln (1 + e^{-t_n w^T x_n}) \quad \text{①}$$

$t_n \in \{-1, 1\}$

when t_n 's are {-1, 1} in sklearn

$$E(w) = \frac{1}{2} w^T w + C \sum_{n=1}^N \ln (1 + e^{-t_n w^T x_n}) \quad \text{②}$$

To get maximum likelihood estimate of parameter w ,

$$\hat{w} = \underset{w}{\operatorname{argmin}} \left[\frac{1}{2} w^T w + \sum_{n=1}^N \ln (1 + e^{-t_n w^T x_n}) \right]$$

$$\Rightarrow \hat{w} = \underset{w}{\operatorname{argmin}} \left[\frac{1}{2} w^T w + \frac{1}{2} \sum_{n=1}^N \ln (1 + e^{-t_n w^T x_n}) \right] \quad (\text{general form})$$

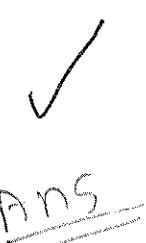
$$\text{but } \hat{w} = \underset{w}{\operatorname{argmin}} \left[\frac{1}{2} w^T w + C \sum_{n=1}^N \ln (1 + e^{-t_n w^T x_n}) \right] \quad (\text{scikit form})$$

(b)

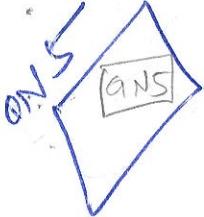
∴

$$C = \frac{1}{2}$$

Ans

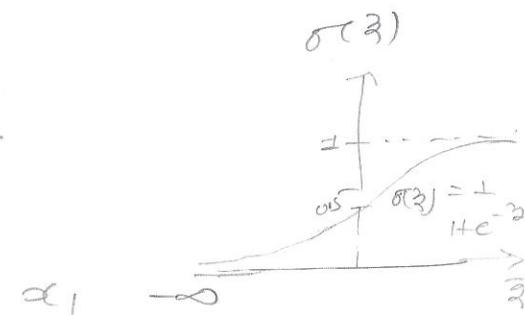
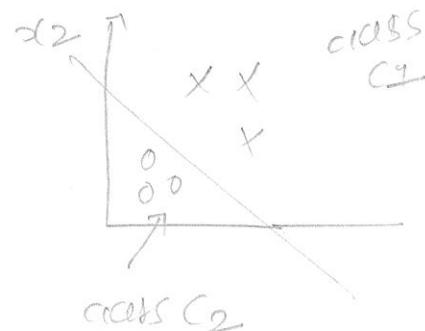


Example
To add bias term
sklearn uses
 $w^T \rightarrow w^T - C \cdot \mathbf{1}$
 E_D



softmax regression is special
case of logistic regression

FOR LR



(sigmoid fn)

(binary classification)

the posterior prob of class C1 is
posterior prob

$$p(c_1|x; w) = p(c_1|x) = \sigma(w^T x) = \frac{1}{1 + e^{-w^T x}}$$

the posterior prob of class C2 is,

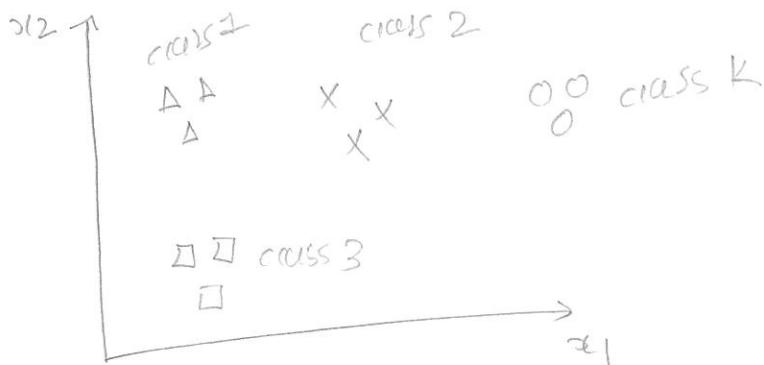
$$p(c_2|x) = 1 - \sigma(w^T x) = \frac{e^{-w^T x}}{1 + e^{-w^T x}}$$

hypothesis → the hypothesis for LR,

$$h(w) = \begin{cases} \frac{1}{1 + e^{-w^T x}} & \rightarrow \text{for class 1} \\ \frac{e^{-w^T x}}{1 + e^{-w^T x}} & \rightarrow \text{for class 2} \end{cases}$$

Again, FOR SR

(softmax Regression
or multiclass logistic regression)



total number of
samples = N

In softmax regression,
the posterior probability of
any class C_K is given by,

$$P(C_K|x, w) = \frac{e^{w_K^T x}}{\sum_{k=1}^K e^{w_k^T x}}$$

posterior
prob.

INPUT DATA X	target
x_1 some values	triangle +1
x_2	cross +2
:	rectangle +3
x_N	circle +N

[Bishop, 4.3.4]
eq 4.1104
p. 209 / Pdtb 227

note $a_K = w_K^T x$ is
called activation of
for the parameter
vector w_K

$$h = \frac{1}{\sum_{k=1}^K e^{w_k^T x}} \begin{bmatrix} e^{w_1^T x} \\ e^{w_2^T x} \\ \vdots \\ e^{w_K^T x} \end{bmatrix}$$

hyp:
plus sign
 $x w^T x$

sigmoid is -
softmax is +

Show $h(w) = h(w - \psi)$

[SL has
overparametrization
property]

$$h(w - \psi) = \frac{e^{(w_K - \psi)^T \alpha}}{\sum_{k=1}^K e^{(w_k - \psi)^T \alpha}}$$

$$= \frac{e^{w_K^T \alpha} \cdot e^{-\psi^T \alpha}}{\sum_k e^{w_k^T \alpha} \cdot e^{-\psi^T \alpha}}$$

$$= \frac{e^{w_K^T \alpha} \cdot e^{-\psi^T \alpha}}{e^{-\psi^T \alpha} \sum_k e^{w_k^T \alpha}}$$

~~$e^{-\psi^T \alpha}$~~
 \sim
does not depend on K

$$= \frac{e^{w_K^T \alpha}}{\sum_k e^{w_k^T \alpha}}$$

$$\boxed{h(w - \psi) = h(w)}$$

∴ If we change any parameter vector $w_K \rightarrow w_K - \psi$
we get same hypothesis for softmax regression.

for two-class SR,

$$h = \frac{1}{\sum_k e^{w_k^T x}} \begin{bmatrix} e^{w_1^T x} \\ e^{w_2^T x} \end{bmatrix}$$

$$= \frac{1}{e^{\vec{w}_1^T x} + e^{\vec{w}_2^T x}} \begin{bmatrix} e^{\vec{w}_1^T x} \\ e^{\vec{w}_2^T x} \end{bmatrix}$$

use overparametrization property of SR

$$h(\vec{w} - \vec{\psi}) = h(\vec{w}) \quad \text{where } \vec{\psi} = \vec{w}_2$$

$$= \frac{1}{e^{(\vec{w}_1 - \vec{w}_2)^T x} + 1} \begin{bmatrix} e^{(\vec{w}_1 - \vec{w}_2)^T x} \\ 1 \end{bmatrix}$$

write $\vec{w}_2 - \vec{w}_1 = \vec{w}$

$$= \frac{1}{1 + e^{-\vec{w}^T x}} \begin{bmatrix} e^{-\vec{w}^T x} \\ 1 \end{bmatrix}$$

$$(h)_{SR} = \begin{bmatrix} \frac{e^{-\vec{w}^T x}}{1 + e^{-\vec{w}^T x}} \\ \frac{1}{1 + e^{-\vec{w}^T x}} \end{bmatrix} = (h)_{LR} \quad \#$$

proved

This is same as h for LR

Direct method to prove softmax ($K=2$) = LR

the hypothesis for softmax regression is

$$h(w, x) = \frac{1}{\sum_{k=1}^K e^{w_k^T x}}$$

$$\begin{bmatrix} e^{w_1^T x} \\ e^{w_2^T x} \\ \vdots \\ e^{w_K^T x} \end{bmatrix}$$

when $K=2$,

$$h = \frac{1}{e^{w_1^T x} + e^{w_2^T x}} \begin{bmatrix} e^{w_1^T x} \\ e^{w_2^T x} \end{bmatrix}$$

$$(h) = \text{softmax} \begin{bmatrix} \frac{e^{w_1^T x}}{e^{w_1^T x} + e^{w_2^T x}} \times e^{-w_2^T x} \\ \frac{e^{w_2^T x}}{e^{w_1^T x} + e^{w_2^T x}} \times e^{-w_1^T x} \end{bmatrix}$$

(20)

✓

$$= \begin{bmatrix} e^{-(w_2^T - w_1^T)x} \\ \frac{e}{1 + e^{-(w_2^T - w_1^T)x}} \\ \frac{1}{1 + e^{-(w_2^T - w_1^T)x}} \end{bmatrix} = \begin{bmatrix} e^{-w^T x} \\ \frac{1}{1 + e^{-w^T x}} \\ \frac{1}{1 + e^{-w^T x}} \end{bmatrix}$$

where
 $w^T = w_2^T - w_1^T$

or,

$$(h)_{\text{softmax}} = \begin{bmatrix} \frac{e^{-w^T x}}{1 + e^{-w^T x}} \\ \vdots \\ \frac{1}{1 + e^{-w^T x}} \end{bmatrix}$$

This is (h) Logistic Regression

$\therefore h_{\text{SR}} = (h)_{\text{LR}}$ proved!

on 6
QNG

derivative of costⁿ for softmax Reg &

for softmax Regression,

prob of data x belonging to class k i.e.

posterior prob of class k given data x is,

posterior
prob
of
class
 k

$$p(k|x_n, w) = \frac{e^{w_k^T x_n}}{\sum_{k=1}^K e^{w_k^T x_n}} \quad (1)$$

BISHOP 4.3.4
p. 209 / 227

rename class k by target t_n , then the probability that the n th example x_n belongs to target t_n is given by,

$$p(t_n|x_n, w) = \frac{e^{w_{t_n}^T x_n}}{\sum_{n=1}^N e^{w_{t_n}^T x_n}} \quad (2)$$

then,
likelihood

$$\ell(w) = \prod_{n=1}^N p(t_n|x_n, w) \quad (2)$$

The -ve log likelihood is,

$$-\ln \ell(w) = -\ln \prod_{n=1}^N \left(\frac{e^{w_{t_n}^T x_n}}{\sum_{n=1}^N e^{w_{t_n}^T x_n}} \right)$$

cost
for
data

$$E_0 = \frac{1}{N} \cdot (-\ln \ell(w)) = -\frac{1}{N} \sum_{n=1}^N \ln \left(\frac{e^{w_{t_n}^T x_n}}{\sum_{n=1}^N e^{w_{t_n}^T x_n}} \right) \quad (3)$$

(Weight
vector
for
each
example)

here, t_n is $t_1, t_2, t_3, \dots, t_N$ (single label)
 x_n is $x_{1,n}, x_{2,n}, x_{3,n}, \dots, x_{N,n}$ (x is row vector)
 w is w_1, w_2, \dots, w_N (w is column)
 (label is row w is column) should be single index

here, instead of $\vec{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_K \end{bmatrix}$ we can group

the data so that there are only K weight vectors $w = [w_1 \ w_2 \ \dots \ w_K]$, then,

$$p(c_k | x) = \frac{e^{w_k^T x}}{\sum_{k=1}^K e^{w_k^T x}} \quad (1)$$

C regrouping weight vectors for each K classes, we can rewrite cost func

and,

$$E_D(w) = -\frac{1}{N} \sum_{n=1}^N \left[\sum_{k=1}^K \delta_{k(t_n)} \right] \ln \left(\frac{e^{w_k^T x_n}}{\sum_{k=1}^K e^{w_k^T x_n}} \right)$$

insert this!

$$= -\frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K \delta_{k(t_n)} \left[w_k^T x_n - \ln \sum_{k=1}^K e^{w_k^T x_n} \right]$$

$$E_D = -\frac{1}{N} \sum_{n=1}^N \left[\sum_{k=1}^K \delta_{k(t_n)} w_k^T x_n - \sum_{k=1}^K \delta_{k(t_n)} \ln \left(\sum_{k=1}^K e^{w_k^T x_n} \right) \right]$$

$$\frac{\partial E_D}{\partial w_j} = -\frac{1}{N} \sum_{n=1}^N \left[\frac{\partial}{\partial w_j} \sum_{k=1}^K \delta_{k(t_n)} w_k^T x_n - \sum_{k=1}^K \delta_{k(t_n)} \ln \left(\sum_{k=1}^K e^{w_k^T x_n} \right) \right]$$

here w_j operates only when $w_j = w_k$ and summation vanishes

$$\frac{\partial E_D}{\partial w_j} = -\frac{1}{N} \sum_{n=1}^N \left[\delta_j(t_n) \alpha_n - \delta_j(t_n) \cdot \frac{1 \cdot e^{w_j^T \alpha_n}}{\sum_{k=1}^K e^{w_k^T \alpha_n}} \cdot \delta_j(t_n) \alpha_n \right]$$

$$\frac{\partial E_D}{\partial w_K} = -\frac{1}{N} \sum_{n=1}^N \left[\delta_{1K}(t_n) \alpha_n - \alpha_n \cdot \delta_{1K}(t_n) \frac{e^{w_K^T \alpha_n}}{\sum_{k=1}^K e^{w_k^T \alpha_n}} \right]$$

$\delta_{K(t_n)} \cdot \delta_{1K(t_n)}$
 $p(c_k | \alpha_n)$ (eq 1)
 we change dummy index from $j \rightarrow K$
 $p(c_k)$ is prob only when label is K
 so basically $\delta_{K(t_n)}$ is 1 for $p(c_k)$

(2) ✓

$$= -\frac{1}{N} \sum_{n=1}^N \left[\delta_K(t_n) \alpha_n - \alpha_n p(c_K | \alpha_n) \right]$$

gradient of E_D

$$\boxed{\nabla_w E_D = -\frac{1}{N} \sum_{n=1}^N (\delta_K(t_n) - p(c_K | \alpha_n)) \alpha_n}$$

(A)

Again $\nabla_w E_w(w) = \frac{\partial}{\partial w} \frac{1}{2} w_K^T w_K = \frac{1}{2} \cdot 2w_K^T = \alpha w_K^T$

$$\boxed{\nabla_w E_w = \alpha w_K^T}$$

(B)

~~Ans~~

$$\nabla_w E = \nabla_w E_D + \nabla_w E_w = -\frac{1}{N} \sum_{n=1}^N (\delta_{1K}(t_n) - p(c_K | \alpha_n)) \alpha_n + \alpha w_K^T$$

$$\text{method 2} \quad E_D(w) = -\frac{1}{N} \sum_n (\frac{e^{w^T x_n}}{\sum_k e^{w^T x_k}}) \quad (\text{grouping n-means})$$

$$E_D(w) = -\frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K \delta_{k(t_n)} \ln \left(\frac{e^{w^T x_n}}{\sum_{k=1}^K e^{w^T x_n}} \right) \quad (\text{grouping K-means})$$

$$\text{let, } E_n = \sum_{k=1}^K \delta_{k(t_n)} \ln \left(\frac{e^{w^T x_n}}{\sum_k e^{w^T x_n}} \right)$$

$$E_n = \sum_k \delta_{k(t_n)} w_i^T x_n - \sum_k \delta_{k(t_n)} \ln \left(\sum_k e^{w^T x_n} \right)$$

$$\frac{\partial E_n}{\partial w_j} = \delta_j(t_n) x_n - \delta_j(t_n) \cdot \frac{1}{\sum_j e^{w_j^T x_n}} \cdot e^{w_j^T x_n} \cdot \delta_j(t_n) x_n$$

$$= \delta_j(t_n) x_n - \frac{e^{w_j^T x_n}}{\sum_j e^{w_j^T x_n}} x_n$$

$$\frac{\partial E_n}{\partial w_k} = \delta_{k(t_n)} x_n - \frac{e^{w_k^T x_n}}{\sum_k e^{w_k^T x_n}} x_n$$

$$\boxed{\frac{\partial E_n}{\partial w_k} = (\delta_{k(t_n)} - p(c_k | x_n)) x_n}$$

$$\therefore \nabla_{w_k} E_D(w) = -\frac{1}{N} \sum_n [\delta_{k(t_n)} - p(c_k | x_n)] x_n \quad @$$

also,

$$\boxed{E_w(w) = \frac{\alpha}{2} w_k^T w_k}$$

$$\nabla_{w_k} E_w(w) = \frac{\partial}{\partial w_k} \frac{\alpha}{2} w_k^T w_k = \alpha w_k$$

$$\therefore \nabla E = -\frac{1}{N} \sum_{n=1}^N (\delta_{k(t_n)} - p(c_k | x_n)) x_n + \alpha w_k$$

$$\boxed{\nabla_w E = \alpha \vec{w}_k - \frac{1}{N} \sum_n (\delta_{k(t_n)} - p(c_k | x_n)) \vec{x}_n}$$

Ans

HWS

Bhishan Paudel

142
150

HW Assignment 5 (Due by 10:30am on Nov 2)

1 Theory (110 points)

1. [Properties of Linear Discriminants, 20 points]

We have proven in class that the distance between origin and the decision hyperplane $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = 0$ is equal with $-w_0 / \|\mathbf{w}\|$. Prove that the margin between a point \mathbf{x} and the same decision hyperplane is equal with $h(\mathbf{x}) / \|\mathbf{w}\|$.

2. [Bonus, 20 points]

Prove the two properties above for the general n -dimensional case.

3. [Fisher Criterion and Least Squares, 30 points]

Show that the Fisher criterion can be written in the vectorized form shown below:

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2} = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

4. [Fisher Criterion (*), 20 points]

Reference the PRML Chapter 4 material available on Blackboard under Content.

Using the definitions of the between-class and within-class covariance matrices given by (4.27) and (4.28), respectively, together with (4.34) and (4.36) and the choice of target values described in Section 4.1.5, show that the expression (4.33) that minimizes the sum-of-squares error function can be written in the form (4.37).

5. [Perceptrons, 40 points]

Consider a training set that contains the following 8 examples:

\mathbf{x}	x_1	x_2	x_3	$t(x)$
$\mathbf{x}^{(1)}$	0	0	0	+1
$\mathbf{x}^{(2)}$	0	1	0	+1
$\mathbf{x}^{(3)}$	1.5	0	-1.5	+1
$\mathbf{x}^{(4)}$	1.5	1	-1.5	+1
$\mathbf{x}^{(5)}$	1.5	0	0	-1
$\mathbf{x}^{(6)}$	1.5	1	0	-1
$\mathbf{x}^{(7)}$	0	0	-1.5	-1
$\mathbf{x}^{(8)}$	0	1	-1.5	-1

(a) Prove that the perceptron algorithm does not converge on this dataset. Do not forget to include the bias.

(b) Consider a kernel perceptron that uses a polynomial kernel $k(\mathbf{x}, \mathbf{y}) = (1 + \mathbf{x}^T \mathbf{y})^d$. What is the smallest degree d for which the kernel perceptron would converge on this dataset?

6. [Perceptrons, 10 points]

A kernel perceptron for binary classification is run for a number of epochs E on a training dataset containing N examples, resulting in the dual parameters $\alpha_1, \alpha_2, \dots, \alpha_N$. What is the total number of mistakes that are made during training?

7. [Matrix Computations, 10 points]

Let $U \in R_{k \times m}$ and $X \in R_{n \times m}$. Let u_i and x_i be the i -th columns of U and X , respectively, for $1 \leq i \leq m$. Prove that $UX^T = \sum_{i=1}^m u_i x_i^T$.

2 Submission

Turn in a hard copy of your homework report at the beginning of class on the due date. On this theory assignment, **clear and complete explanations and proofs of your results are as important as getting the right answer.**

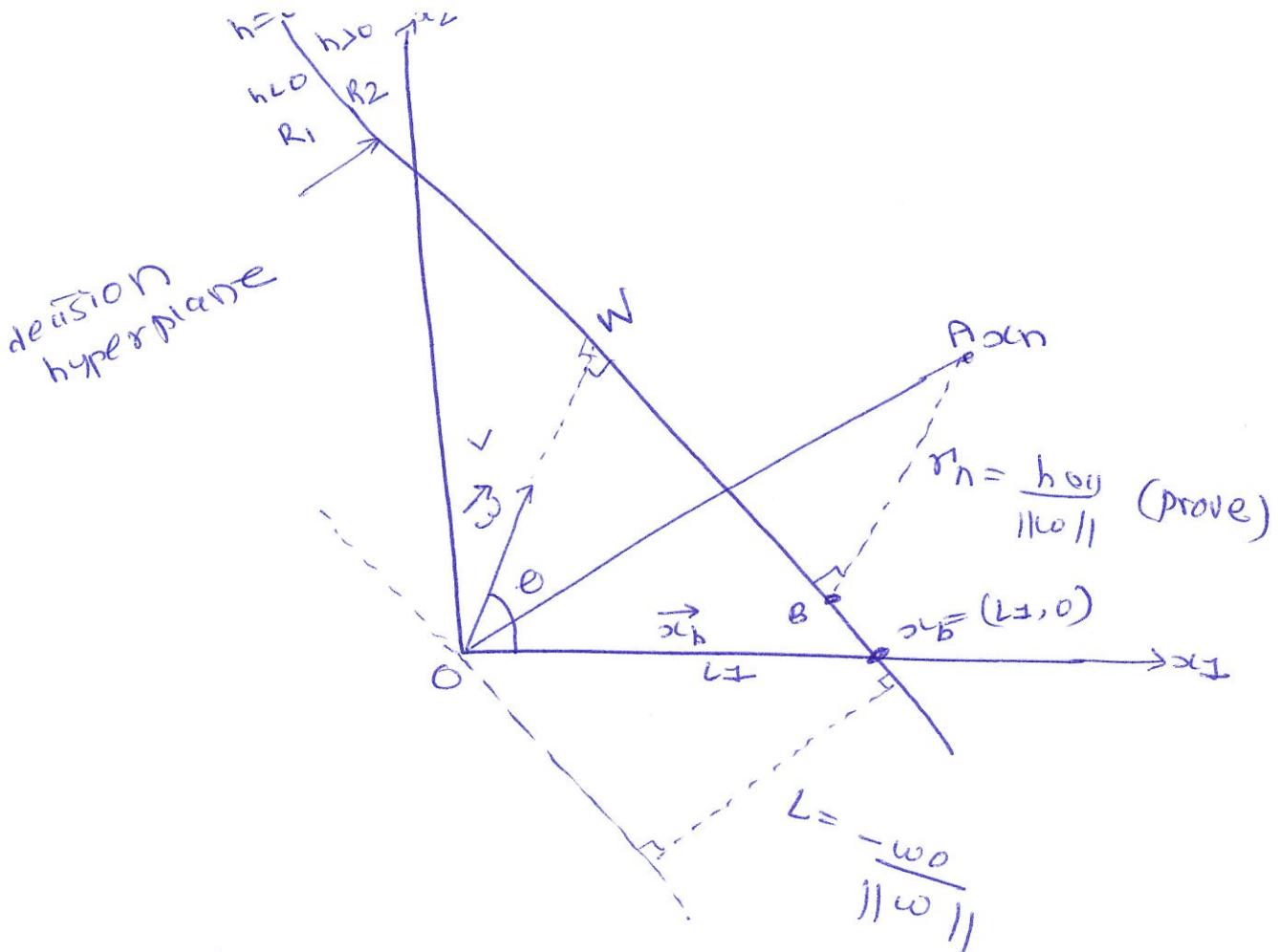
Hw 05

Bhishan Poudel

QNT

prove:

$$\text{dist betw } \underline{\text{data point}} \& \underline{\text{decision hyp.}} = \frac{\|w\|}{\|w\|}$$



$$(a) \text{ move } L_1 = -\frac{w_0}{\|w\|}$$

$$\text{Here, } \vec{w} = [w_1, w_2]$$

$$\vec{x}_b = [x_1, x_2] = [L_1, 0]$$

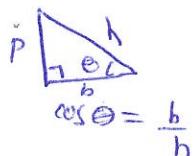
$$\vec{w} \cdot \vec{x}_b + w_0 = 0 = w_1 x_1 + w_2 x_2 + w_0$$

$$0 = w_1 L_1 + w_2 \cdot 0 + w_0$$

$$\boxed{L_1 = -\frac{w_0}{w_1}} \quad \text{--- (1)}$$

again,

$$\cos \theta = \frac{\text{base}}{\text{hypotenuse}} = \frac{L}{\|w\|}$$



$$\cos \theta = \frac{L}{\|w\|}$$

$$\boxed{L = -\frac{w_0 \cos \theta}{w_1}} \quad \text{--- (2)}$$

$$L = -\frac{w_0}{\|w\|} \cdot \frac{\vec{w} \cdot \vec{z}_b}{\|\vec{w}\| \|z_b\|} \quad ; \quad \vec{a} \cdot \vec{b} = \|a\| \|b\| \cos \theta$$

$$= -\frac{w_0}{\|w\|} \cdot \frac{\cancel{(w_1 w_2) + (\cancel{w_1}) \cancel{w_2})}}{\|\vec{w}\| \|z_b\|}$$

$$= -\frac{w_0}{\|w\|} \frac{w_1 k_1}{\|\vec{w}\| \|z_b\|}$$

$L = -\frac{w_0}{\|\vec{w}\|}$

Ans

Again, To find the coordinate of point B,

$$\vec{OB} = \vec{OA} - \vec{AB}$$

$$\vec{B} = \vec{x}_n - r_n \vec{w} \quad \rightarrow \vec{w} \text{ is unit vector perpendicular to decision hyperplane}$$

$B = x_n - r_n \frac{w}{\|w\|}$

But, point B lies in decision boundary, so,

$$\vec{w} \cdot \vec{B} + w_0 = 0$$

$$w^T (x_n - r_n \frac{w}{\|w\|}) + w_0 = 0$$

$$\text{or, } w^T x_n - r_n \frac{w^T w}{\|w\|} + w_0 = 0$$

$$\text{or, } w^T x_n - r_n / \|w\| + w_0 = 0$$

$$\therefore w^T w = \|w\|^2$$

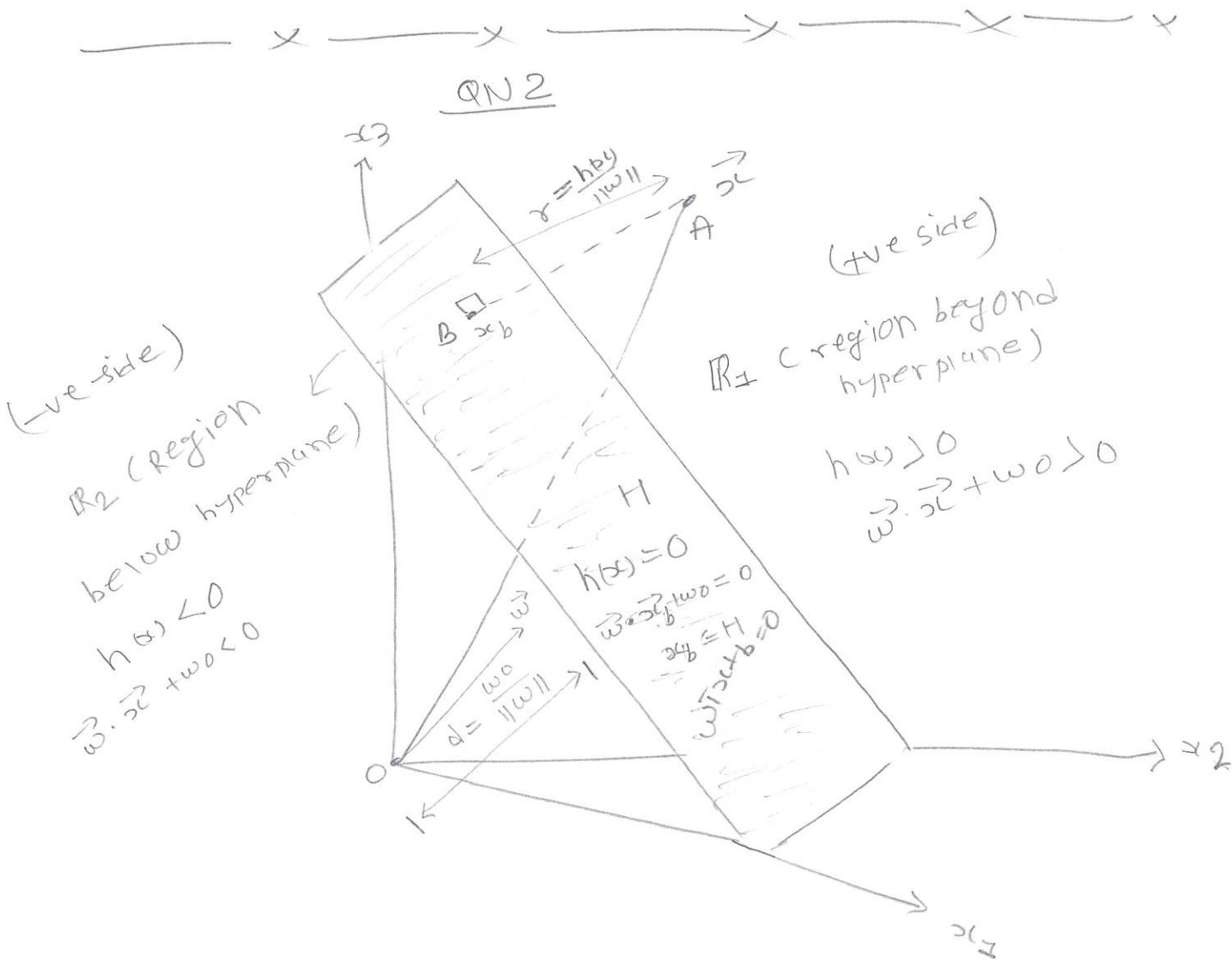
$$\text{or, } \vec{w}^T \vec{x}_n + w_0 = \gamma_n \|\vec{w}\|$$

$$\therefore \gamma_n = AB = \frac{\vec{w}^T \vec{x}_n + w_0}{\|\vec{w}\|}$$

where $h_{\vec{w}} = \vec{w}^T \vec{x} + w_0$

$$\boxed{\gamma_n = \frac{h_{\vec{w}}}{\|\vec{w}\|}} \quad \text{proved!}$$

✓ 20



QN2

For n-dimensional case,
prove:

$$\text{dist betn origin and decision hyperplane} = -\frac{w_0}{\|w\|}$$

$$\& \text{dist betn datapoint and decision hyperplane} = \frac{h(x)}{\|w\|}$$

Soln Let \vec{x} be any point on the positive side
of ^{decision} boundary (hyperplane), then,

$$\vec{x} = \vec{x}_b + r \frac{\vec{w}}{\|w\|} \quad \text{where } \frac{\vec{w}}{\|w\|} \text{ is unit vector}$$

$$\boxed{\vec{x}_b = \vec{x} - r \frac{\vec{w}}{\|w\|}} \quad \text{①}$$

perpendicular to decision
hyperplane

but, \vec{x}_b lies in boundary, we need to find r ,

$$\boxed{\vec{w} \cdot \vec{x}_b + w_0 = 0} \quad \text{②}$$

✓

$$\text{or, } w^T (\vec{x} - r \frac{\vec{w}}{\|w\|}) + w_0 = 0 = w^T \vec{x} - r \frac{w^T w}{\|w\|} + w_0 = 0$$

$$\text{if, } w^T \vec{x} + w_0 = r \|w\| = h(x)$$

✓

distance from
datapoint
to
hyperplane

$$\Rightarrow r = \frac{w^T \vec{x} + w_0}{\|w\|}$$

$$\boxed{r = \frac{h(x)}{\|w\|}}$$

DNS

Also, we have to show the distance of hyperplane from the origin,

$$\vec{OB} = -\frac{\underline{w}_0}{\|w\|} \vec{OB}$$

Here, the point x_b lies in decision hyperplane, so,

$$\vec{w} \cdot \vec{x}_b + w_0 = 0$$

$$w \cdot x_b + w_0 = 0$$

$$w \cdot x_b = -w_0$$

distance from
origin to
decision hyperplane

$$x_b = -\frac{\underline{w}_0}{\|w\|}$$

Ans

QN3

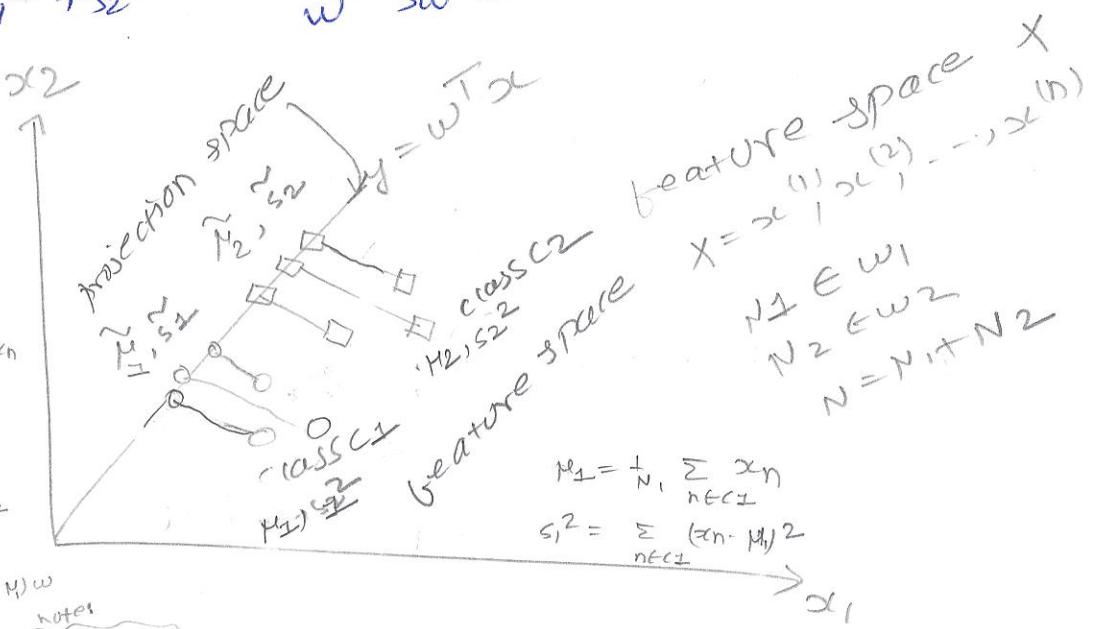
show that Fisher criterion

Bishop
(ex 4.15) can be written as,

$$J(w) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2} = \frac{w^T S_B w}{w^T S_w w}$$

Solution

$$\begin{aligned}\hat{\mu}_1 &= \frac{1}{N_1} \sum_{n \in C_1} w^T x_n \\ &= w^T m_1 \\ s_1^2 &= \sum_{n \in C_1} (y_n - \hat{\mu}_1)^2 \\ &= \sum_{n \in C_1} (w^T x_n - w^T \mu_1)^2 \\ &= \sum_{n \in C_1} w^T (x_n - \mu_1)^\top (x_n - \mu_1) w \\ \hat{\mu}_2 &= w^T \mu_2 \\ (\text{big } S)\end{aligned}$$



In feature space X , we have input data,

$$x = (x^{(1)}, x^{(2)}, \dots, x^{(n)})$$

feature space X , mean

$$\mu_i = \frac{1}{N_i} \sum_{x \in C_i} x$$

projection space Y , mean $\tilde{\mu}_i = \frac{1}{N_i} \sum_{y \in C_i} y$

$$= \frac{1}{N_i} \sum_{x \in C_i} w^T x$$

$$\tilde{\mu}_i = w^T \mu_i$$

Also,

Variance,
(feature space)

Variance
(projection
space)

$$\sigma_i^2 = \frac{1}{N_i} \sum_{x \in \mathcal{E}_i} (\alpha - \mu_i)^2$$

$$\tilde{\sigma}_i^2 = \frac{1}{N_i} \sum_{y \in \mathcal{E}_i} (y - \tilde{\mu}_i)^2$$

Variance

Also, scatter in feature space,

$$S_i^2 = \sum_{x \in \mathcal{E}_i} (\alpha - \mu_i)^2 = \sum_{x \in \mathcal{E}_i} (\alpha - \mu_i)(\alpha - \mu_i)^T$$

Scatter in projection space,

$$\tilde{S}_i^2 = \sum_{y \in \mathcal{E}_i} (y - \tilde{\mu}_i)^2$$

Scatter

Also,

within-class scatter

$$S_w = S_1 + S_2$$

Also,

$$\tilde{S}_w = \tilde{S}_1 + \tilde{S}_2$$

Now, we define Fisher's criterion,

$$J(w) = \frac{(\tilde{\mu}_1 - \tilde{\mu}_2)^2}{\tilde{s}_1^2 + \tilde{s}_2^2}$$

Fisher's
linear
discriminant

NOW,
 $(\tilde{\mu}_1 - \tilde{\mu}_2)^2 = (\bar{w}^T \mu_1 - \bar{w}^T \mu_2)^2$
 $= \bar{w}^T (\mu_1 - \mu_2) (\mu_1 - \mu_2)^T w$

$$(\tilde{\mu}_1 - \tilde{\mu}_2)^2 = \bar{w}^T S_B w \quad \textcircled{A}$$

where, $S_B = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T$

is between-class scatter.

Also, $\tilde{s}_i^2 = \sum_{y \in \mathcal{C}_i} (y - \tilde{\mu}_i)^2$
 $= \sum_{x \in \mathcal{C}_i} (\bar{w}^T x - \bar{w}^T \mu_i)^2$

$$= \sum_{x \in \mathcal{C}_i} \bar{w}^T (x - \mu_i) (\bar{x} - \bar{\mu}_i)^T w$$

$$\tilde{s}_i^2 = \bar{w}^T S_i w$$

where, feature-space scatter

$$S_i = \sum_{x \in \mathcal{C}_i} (x - \mu_i)(x - \mu_i)^T$$

$$\Rightarrow \tilde{s}_1^2 = \bar{w}^T S_1 w$$

$$\tilde{s}_2^2 = \bar{w}^T S_2 w$$

$$\tilde{s}_1^2 + \tilde{s}_2^2 = \bar{w}^T (S_1 + S_2) w = \bar{w}^T S_w w$$

where $S_w = S_1 + S_2$ is within-class scatter.

$$\Rightarrow \boxed{\xi_1^2 + \xi_2^2 = \omega^T S_w \omega} \quad \text{②}$$

Using equations ① & ②, we have,

$$J(\omega) = \frac{|\tilde{\mu}_1 - \tilde{\mu}_2|^2}{\xi_1^2 + \xi_2^2}$$

$$\boxed{J(\omega) = \frac{\omega^T S_B \omega}{\omega^T S_w \omega}}$$

Q.E.D.

30

(Q4)

In Bishop Book, (P. 190 or 208)

Bishop use equations: 4.27, 4.28, 4.34, 4.36, 4.33
 ex 4.6
 P.

Derive: $\left(S_w + \frac{N_1 N_2}{N} S_B \right) w = N (m_1 - m_2)$

4.37

Here, given, SSE is given by

$$E(w) = \frac{1}{2} \sum_{n=1}^N (w^T x_n + w_0 - t_n)^2$$

— (4.31) P. 190/208

$$\frac{\partial E}{\partial w} = 0 = \sum_{n=1}^N (w^T x_n + w_0 - t_n) x_n \quad — (4.33)$$

(4.34)
 $w_0 = -w^T m$

$$t_n = \begin{cases} \frac{N_1}{N} & \text{if } x_n \in C_1 \\ -\frac{N_2}{N} & \text{if } x_n \in C_2 \end{cases}$$

so that $\sum_{n=1}^N t_n = N_1 \frac{N}{N} - N_2 \frac{N}{N}$

mean of the total dataset (population mean)

(4.36)
 $m = \frac{1}{N} \sum_{n=1}^N x_n = \frac{1}{N} N_1 m_1 + N_2 m_2$

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

where, $m_1 = \frac{1}{N_1} \sum_{n \in C_1} x_n$ (sample mean of class 1)
 (mean of x_n belonging to class C_1)

(4.22)
 P. 205
 $m_2 = \frac{1}{N_2} \sum_{n \in C_2} x_n$

$$\Rightarrow \sum_{n \in C_1} x_n = N_1 m_1$$

$$\Rightarrow \sum_{n \in C_2} x_n = N_2 m_2$$

we define within-class scatter as,

$$\text{eq 4.28} \quad S_w = \sum_{n \in C_1} (x_n - m_1)(x_n - m_1)^T + \sum_{n \in C_2} (x_n - m_2)(x_n - m_2)^T$$

and, between-class scatter as,

$$\text{eq 4.29} \quad S_B = (m_2 - m_1)(m_2 - m_1)^T$$

we have to prove,

~~to prove~~
$$(S_w + \frac{N_1 N_2}{N} S_B)w = N(m_1 - m_2) \quad (\text{eq 4.37})$$

$$\Rightarrow \frac{N(m_1 - m_2)}{w} = S_w + \frac{N_1 N_2 S_B}{N} \quad \square$$

we start from gradient zero equation (4.33),

$$0 = \sum_{n=1}^N (w^T x_n - w^T m - t_n) x_n$$

$$= \sum_{n=1}^N (w^T x_n - w^T m - t_n) x_n$$

$$= \sum_{n=1}^N w^T (x_n x_n^T - m x_n) - \sum_{n=1}^N t_n x_n$$

$$= \sum_{n \in C_1} (x_n x_n^T - x_n m^T) w - \sum_{n \in C_1} t_n x_n \quad \begin{array}{l} \uparrow \\ \text{use } t_n = \begin{cases} N_1 & \text{for } C_1 \\ N_2 & \text{for } C_2 \end{cases} \end{array}$$

$$+ \sum_{n \in C_2} (x_n x_n^T - x_n m^T) w - \sum_{n \in C_2} t_n x_n$$

$$m_2 = \frac{1}{N_2} \sum_{n \in C_2} x_n \Rightarrow \sum_{n \in C_2} x_n = N_2 m_2$$

$$0 = \left(\sum_{n \in C_1} x_n x_n^T - N_1 m_1 m_1^T \right) w - \frac{N_1}{N} m_1$$

$$\left(\sum_{n \in C_2} x_n x_n^T - N_2 m_2 m_2^T \right) w + \frac{N_2}{N} m_2$$

$$0 = \left(\sum_{n \in C_1} x_n x_n^T + \sum_{n \in C_2} x_n x_n^T - (N_1 m_1 + N_2 m_2) m^T \right) w - N(m_1 - m_2)$$

$$\frac{N(m_1 - m_2)}{w} = \underbrace{\sum_{n \in C_1} x_n x_n^T + \sum_{n \in C_2} x_n x_n^T}_{\text{write in terms of } S_w} - (N_1 m_1 + N_2 m_2) m^T$$

②

$$m = \frac{N_1 m_1 + N_2 m_2}{N}$$

NOW, we expand the first term in S_w ,

$$\begin{aligned} \sum_{n \in C_1} (x_n - m_1) (x_n - m_1)^T &= \sum_{n \in C_1} (x_n x_n^T - x_n m_1^T - m_1 x_n^T + m_1 m_1^T) \\ &= \sum_{n \in C_1} x_n x_n^T - m_1^T \sum_{n \in C_1} x_n - m_1 \sum_{n \in C_1} x_n^T + \sum_{n \in C_1} m_1 m_1^T \\ &= \sum_{n \in C_1} x_n x_n^T - N_1 m_1 m_1^T - N_1 m_1 m_1^T + N_1 m_1 m_1^T \end{aligned}$$

$$\boxed{\sum_{n \in C_1} (x_n - m_1) (x_n - m_1)^T = \sum_{n \in C_1} x_n x_n^T - N_1 m_1 m_1^T}$$

$$\text{So, } S_w = \sum_{n \in C_1} (x_n - m_1) (x_n - m_1)^T + \sum_{n \in C_2} (x_n - m_2) (x_n - m_2)^T$$

$$\boxed{S_w = \sum_{n \in C_1} x_n x_n^T + \sum_{n \in C_2} x_n x_n^T - N_1 m_1 m_1^T - N_2 m_2 m_2^T}$$

$$\Rightarrow \sum_{n \in C_1} x_n x_n^T + \sum_{n \in C_2} x_n x_n^T = S_w + N_1 m_1 m_1^T + N_2 m_2 m_2^T$$

then eqn ② becomes

$$\frac{N(m_1 - m_2)}{w} = sw + \underbrace{N_1 m_1 m_1^T}_{\text{1}} + \underbrace{N_2 m_2 m_2^T}_{\text{2}} - (N_1 m_1 + N_2 m_2) \cdot \underbrace{(N_1 m_1 + N_2 m_2)}_{N}^T$$

$$= sw + m_1 m_1^T \left(N_1 - \frac{N_1^2}{N} \right) + m_2 m_2^T \left(N_2 - \frac{N_2^2}{N} \right)$$

$$- \frac{N_1 N_2}{N} (m_1 m_2^T + m_2 m_1^T)$$

$$N = N_1 + N_2$$

$$= sw + \frac{(N_1 + N_2)N_1 - N_1^2}{N} m_1 m_1^T + \frac{(N_1 + N_2)N_2 - N_2^2}{N} m_2 m_2^T$$

$$- \frac{N_1 N_2}{N} (m_1 m_2^T + m_2 m_1^T)$$

$$= sw + \frac{N_1 N_2}{N} (m_1 m_1^T + m_2 m_2^T)$$

$$- \frac{N_1 N_2}{N} (m_1 m_2^T + m_2 m_1^T)$$

$$= sw + \frac{N_1 N_2}{N} (m_1 m_1^T + m_2 m_2^T - m_1 m_2^T - m_2 m_1^T)$$

$$\frac{N(m_1 - m_2)}{w} = sw + \frac{N_1 N_2}{N} s_B$$

$$\begin{aligned} s_B &= (m_2 - m_1)(m_2 - m_1)^T \\ &= m_2 m_2^T - m_2 m_1^T - m_1 m_2^T + m_1 m_1^T \\ &= m_1 m_1^T + m_2 m_2^T - m_1 m_2^T - m_2 m_1^T \end{aligned}$$

$$\therefore \boxed{\left(sw + \frac{N_1 N_2}{N} s_B \right) w = N(m_1 - m_2)}$$

Q.E.D.

✓

20

QNS

Given training set,

x	x_1	x_2	x_3	$t(x)$
$x^{(1)}$	0	0	0	+1
$x^{(2)}$	0	1	0	+1
$x^{(3)}$	1.5	0	-1.5	+1
$x^{(4)}$	1.5	1	-1.5	+1
$\bar{x}^{(5)}$	1.5	0	0	-1
$x^{(6)}$	1.5	1	0	-1
$x^{(7)}$	0	0	-1.5	-1
$x^{(8)}$	0	1	-1.5	-1

a) prove that the perceptron algorithm does not converge on this dataset. (Note: include bias term)

b) consider a kernel perceptron that uses a polynomial kernel $K(x,y) = (1+x^T y)^d$. What is the smallest degree d for which the kernel perceptron would converge on this dataset?

(50)

for a given example, example belongs to class + if

$$w_0 + \vec{w} \cdot \vec{x}_n \geq 0$$

and, example belongs to class - if

$$w_0 + \vec{w} \cdot \vec{x}_n < 0$$

now, $\vec{x}^{(1)} = [w_1 \ w_2 \ w_3 \ 1] \quad t^{(1)} = +1$

$$w_0 + w_1 \cdot 0 + w_2 \cdot 0 + w_3 \cdot 0 \geq 0$$

$$w_0 \geq 0$$

similarly,

w

we {

$$\begin{aligned} ① \quad w_0 + 0 + 0 + 0 &\geq 0 \quad \text{--- (1)} \\ ② \quad w_0 + 0 + w_2 + 0 &\geq 0 \quad \text{--- (2)} \\ ③ \quad w_0 + +\sqrt{w_1} + 0 - +\sqrt{w_3} &\geq 0 \quad \text{--- (3)} \\ ④ \quad w_0 + +\sqrt{w_1} + w_2 - +\sqrt{w_3} &\geq 0 \quad \text{--- (4)} \\ ⑤ \quad w_0 + +\sqrt{w_1} + 0 + 0 &< 0 \quad \text{--- (5)} \\ ⑥ \quad w_0 + +\sqrt{w_1} + w_2 + 0 &< 0 \quad \text{--- (6)} \\ ⑦ \quad w_0 + 0 + 0 - +\sqrt{w_3} &< 0 \quad \text{--- (7)} \\ ⑧ \quad w_0 + 0 + w_2 - +\sqrt{w_3} &< 0 \quad \text{--- (8)} \end{aligned}$$

④ - ③ $\Rightarrow w_2 \geq 0$
 ⑥ - ⑤ $\Rightarrow w_2 \leq 0$

✓

here, ① - ③ gives $w_2 \geq 0$

but ⑥ - ⑤ gives $w_2 \leq 0$

this is
contradictory

\Rightarrow
perceptron
fails.

(5b) Here, the given dataset is not linearly separable in x -space.

We shall propose a perceptron that uses polynomial kernel;

The polynomial kernel is,

$$K(\mathbf{x}_i, \mathbf{x}') = (\mathbf{a}^T \mathbf{x}' + b)^d$$

for $a = b = 1$

$$K(\mathbf{x}_i, \mathbf{x}') = (1 + \mathbf{x}^T \mathbf{x}')^d$$

a = scaling factor

b = bias term

d = degree of polynomial

\mathbf{x} = augmented dataset
of shape N rows
 $m+1$ cols

\mathbf{x}' = one example of dataset
with $m+1$ features

when $d=1$: the kernel is linear classifier, but
our dataset is NON linear, so it can not
classify correctly. 12

When $d=2$: when degree is 2, quadratic kernel
can classify linearly non-separable
dataset. and find a vector to
separate the data set.

when $d \geq 2$: higher degree polynomial kernels can
fit the linearly non-separable dataset
good, but, they may overfit the
test dataset.
may
overfit

∴ smallest degree $d=2$ ANS

Qn 6

For a kernel perceptron,

from lecture note, (lecture 06)

the algorithm is

$$1. f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} = \sum_n \alpha_n t_n \mathbf{x}^T \mathbf{x} = \sum_n \alpha_n t_n K(\mathbf{x}, \mathbf{x})$$

kernel

2. initialized: $\alpha_n = 0$

for e in E epochs

3. for $n=1 \dots N$

$$h_n = \text{sgn}(f(\mathbf{x}_n))$$

if $h_n \neq t_n$ then

$$\alpha_n = \alpha_n + 1$$

Now, from step 3, if we run perceptron

E number of epochs,

then total number of mistakes made

so far if,

$$M = \sum_{n=1}^N \alpha_n$$

Ans

/D



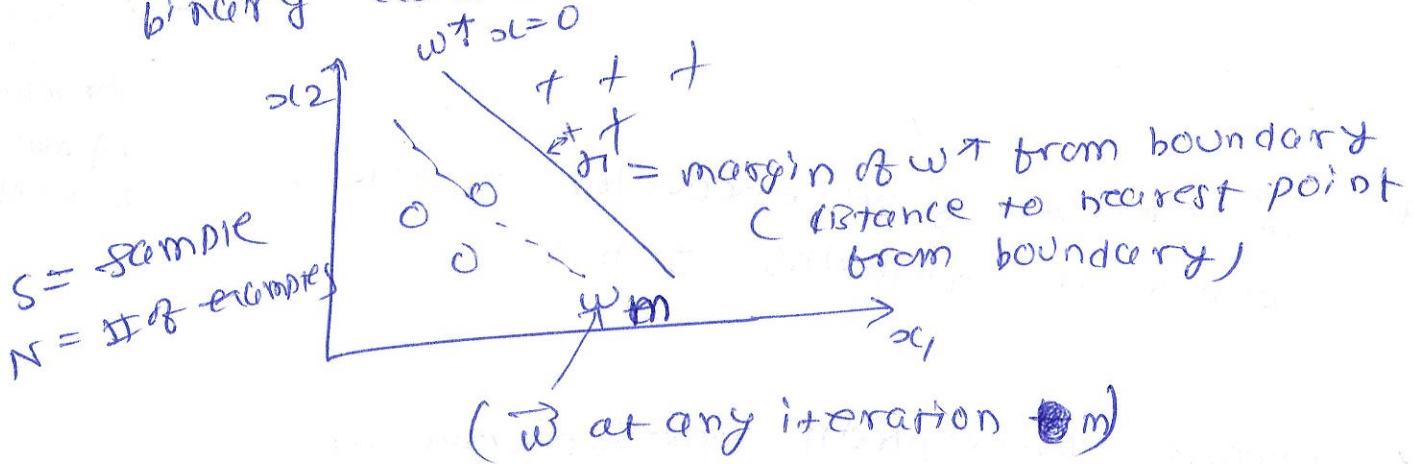
Q6B

Method 2 (Extra practice) A Kernel perceptron for binary classification

- maximum error bound
classification is run for a number of epochs E
on a training dataset containing N examples,
resulting in the dual parameters $\alpha_1, \alpha_2, \dots, \alpha_N$
what is the total number of mistakes that are made during training?

Soln consider a kernel perceptron for

binary classification.



Let there is a unit vector \vec{w}^T that can separate N training examples into two classes with margin γ_1 .

Let $R = \text{norm of maximum value of } \vec{x}^T$ in dataset

$$\|\vec{R}\| = \max_{x \in S} \|\vec{x}\|$$

perceptron algorithm:

1. initialize w to zero. $w_1 = [0 \ 0 \ 0 \ 0]$

2. for e in epochs

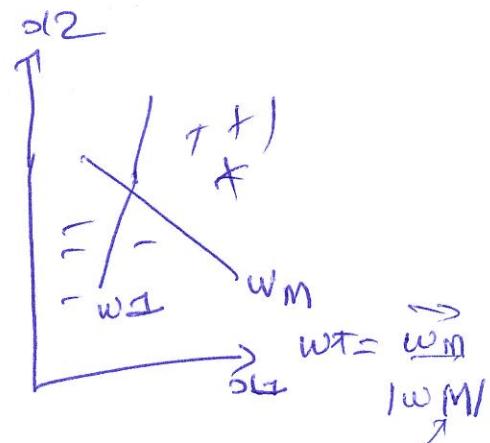
for $n = 1$ to N

if $t_n(w \cdot x_n) < 0$

$$w = w + t_n x_n$$

if converged w :

break



for total
M mistake
 $1 \leq m \leq M$

* \vec{w} changes if we encounter mistake, otherwise remains same.

* suppose perceptron makes a mistake at iteration m , then,

$$\vec{w}_m \cdot \vec{w}^* = (\vec{w}_{m-1} + t \vec{x}) \cdot \vec{w}^*$$

$$= \vec{w}_{m-1} \cdot \vec{w}^* + t (\vec{x} \cdot \vec{w}^*)$$

$$w_m \cdot w^* \geq \vec{w}_{m-1} \cdot \vec{w}^* + \gamma$$

from definition of γ .
each data point x
is $\geq \gamma$ from the
boundary w^*

$$\therefore \boxed{w_m \cdot w^T \geq w_{m-1} \cdot w^T + \gamma} \quad \text{--- (A)}$$

Also,

$$\|w_m\|^2 = \|w_{m-1} + t\vec{x}\|^2$$

$$= \|w_{m-1}\|^2 + 2t(w_{m-1} \cdot \vec{x}) + \|t\vec{x}\|^2$$

④ maximum is a
iteration

⑤ Here we have mistake at iteration

m , so the previous iteration

has hypothesis

$$t(w_{m-1} \cdot \vec{x}) < 0$$

$$\boxed{\|w_m\|^2 \leq \|w_{m-1}\|^2 + R^2} \quad \text{(from ④ & ⑤)}$$

Here, ~~④~~ inequalities ④ & ⑤ hold true
for all the iterations m ,

$$w_m \cdot w^T \geq w_{m-1} \cdot w^T + \gamma$$

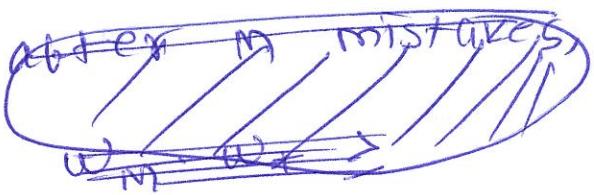
$$\|w_m\|^2 \leq \|w_{m-1}\|^2 + R^2$$

$$\|w_m\|^2 - \|w_{m-1}\|^2 \leq R^2$$

difference Δ

C If there is a mistake then square of weight vectors is less than R^2),

(for each mistake $w^2 - w_{m-1}^2$ cannot be larger than R^2)



for all iterations,

$$\vec{w}_m \cdot \vec{w}^\top \geq w_{m-1} + \gamma$$

after M mistakes,

$$\boxed{\vec{w}_M \cdot \vec{w}^\top \geq r_M} \quad \text{--- } \textcircled{C}$$

also, for all iterations,

$$\|\vec{w}_m\|^2 \leq \|\vec{w}_{m-1}\|^2 + R^2$$

after M mistakes

$$\boxed{\begin{aligned} \|\vec{w}_M\|^2 &\leq R^2 M \\ \|\vec{w}_M\| &\leq R \sqrt{M} \end{aligned}} \quad \text{--- } \textcircled{D}$$

so, from \textcircled{C} & \textcircled{D} ,

$$r_M \leq \vec{w}_M \cdot \vec{w}^\top$$

$$\leq \|\vec{w}_M\| \|\vec{w}^\top\|$$

$$\leq \|\vec{w}_M\|$$

$$r_M \leq R \sqrt{M}$$

$$\pi \sqrt{M} \leq R$$

$$\pi^2 M \leq R^2$$

$$\therefore \vec{a} \cdot \vec{b} = ab \cos \theta$$

~~cosine is less than~~

or equal to 1

\vec{w}^\top is unit vector of \vec{w}_M

$$\|\vec{w}^\top\| = 1$$

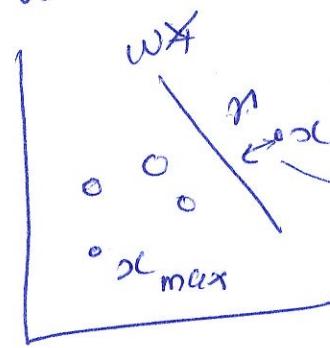
$$\boxed{M \leq \left(\frac{R}{\pi}\right)^2}$$

Ans

Hence,

let the data (training) sample S is,

$$\begin{array}{ll} \alpha_1 & t_1 \\ \alpha_2 & t_2 \\ \alpha_3 & t_3 \\ \vdots & \vdots \\ \alpha_N & t_N \end{array}$$



nearest point
to boundary
 $w^T \cdot \vec{\alpha} = 0$

example,

$$\|\alpha_{\max}\| = R$$

$$\begin{array}{lll} \alpha_1: 0 0 1 2 3 & t_1 = + \\ \alpha_2: 5 0 3 2 6 & t_2 = - \\ \alpha_N: 10 20 30 40 & t_N = + \end{array}$$

then, if we use Perceptron to train our model, the total number of mistakes (M)

made is,

$$M \leq \left(\frac{R}{\gamma_1}\right)^2$$

QNT

Let $V \in R_{k \times m}$

$X \in R^{n \times m}$

$v_i, x_i = i^{\text{th}} \text{ column of } V \text{ and } X \text{ respectively}$
 $1 \leq i \leq m$.

prove $VX^T = \sum_{i=1}^m v_i x_i^T$

SOLUTION:

$$V = V_{k,m} = \begin{pmatrix} v_{11} & v_{12} & \cdots & v_{1m} \\ v_{21} & v_{22} & \cdots & v_{2m} \\ \vdots & & & \\ v_{k1} & v_{k2} & \cdots & v_{km} \end{pmatrix}_{k,m}$$

$$X = X_{n,m}$$

n rows
m cols

$$= \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1m} \\ x_{21} & x_{22} & \cdots & x_{2m} \\ \vdots & & & \\ x_{n1} & x_{n2} & \cdots & x_{nm} \end{pmatrix}_{n,m}$$

\uparrow
~~X~~ column vector

$$V = X^T$$

$$= \begin{pmatrix} x_{11} & x_{21} & \cdots & x_{n1} \\ x_{12} & x_{22} & \cdots & x_{n2} \\ \vdots & \vdots & & \vdots \\ x_{1m} & x_{2m} & \cdots & x_{nm} \end{pmatrix} = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_m^T \end{pmatrix}_{m,1}$$

From definition of matrix multiplication,

$$\text{if } U = U_{km} \\ \text{and } V = V_{mn}$$

then,

$$(UV)_{pq} = \sum_{i=1}^m U_{pi} V_{iq} \quad \text{if } 1 \leq p \leq k \text{ and } 1 \leq q \leq n$$

$1 \leq p \leq k$
 $1 \leq q \leq n$

where $X = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1m} \\ x_{21} & x_{22} & \dots & x_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nm} \end{pmatrix}_{n \times m}$

$n \text{ rows}$
 $m \text{ cols}$

$$x_1 = \begin{pmatrix} x_{11} \\ x_{21} \\ \vdots \\ x_{m1} \end{pmatrix}_{1 \times m}$$

$$= \left(\begin{array}{ccc} u_{11} & u_{12} & u_{1m} \\ u_{21} & u_{22} & u_{2m} \\ \vdots & \vdots & \vdots \\ u_{k1} & u_{k2} & u_{km} \end{array} \right)_{k \times m} \downarrow \left(\begin{array}{ccc} x_{11} & x_{21} & \dots & x_{m1} \\ x_{12} & x_{22} & \dots & x_{m2} \\ \vdots & \vdots & \ddots & \vdots \\ x_{1m} & x_{2m} & \dots & x_{mm} \end{array} \right)_{m \times n}$$

$$= \left(\begin{array}{c} u_{11}x_{11} + u_{12}x_{12} + u_{1m}x_{1m} \\ u_{21}x_{21} + u_{22}x_{22} + u_{2m}x_{2m} \\ \vdots \\ u_{k1}x_{k1} + u_{k2}x_{k2} + u_{km}x_{km} \end{array} \right)_{k \times 1} + \left(\begin{array}{c} u_{11}x_{21} + u_{12}x_{22} + u_{1m}x_{2m} \\ u_{21}x_{21} + u_{22}x_{22} + u_{2m}x_{2m} \\ \vdots \\ u_{k1}x_{k1} + u_{k2}x_{k2} + u_{km}x_{2m} \end{array} \right)_{k \times 1} + \left(\begin{array}{c} u_{11}x_{m1} + u_{12}x_{m2} + u_{1m}x_{mm} \\ u_{21}x_{m1} + u_{22}x_{m2} + u_{2m}x_{mm} \\ \vdots \\ u_{k1}x_{m1} + u_{k2}x_{m2} + u_{km}x_{mm} \end{array} \right)_{k \times 1}$$

$$= \left(\begin{array}{ccc} u_{11}x_{11} & u_{11}x_{21} & u_{11}x_{m1} \\ u_{21}x_{11} & u_{21}x_{21} & u_{21}x_{m1} \\ \vdots & \vdots & \vdots \\ u_{k1}x_{11} & u_{k1}x_{21} & u_{k1}x_{m1} \end{array} \right)_{k \times m} + \left(\begin{array}{ccc} u_{12}x_{12} & u_{12}x_{22} & u_{12}x_{m2} \\ u_{22}x_{12} & u_{22}x_{22} & u_{22}x_{m2} \\ \vdots & \vdots & \vdots \\ u_{k2}x_{12} & u_{k2}x_{22} & u_{k2}x_{m2} \end{array} \right)_{k \times m} + \left(\begin{array}{ccc} u_{1m}x_{1m} & u_{1m}x_{2m} & u_{1m}x_{mm} \\ u_{2m}x_{1m} & u_{2m}x_{2m} & u_{2m}x_{mm} \\ \vdots & \vdots & \vdots \\ u_{km}x_{1m} & u_{km}x_{2m} & u_{km}x_{mm} \end{array} \right)_{k \times m}$$

$$= U_1 X_1^T + U_2 X_2^T + U_m X_m^T$$

10 ✓

similarly,

$$U_1 X_1^T = \begin{pmatrix} U_{11} \\ U_{21} \\ \vdots \\ U_{K1} \end{pmatrix} \begin{pmatrix} X_{11} \\ X_{21} \\ \vdots \\ X_{n1} \end{pmatrix}^T$$

$$= \begin{pmatrix} \rightarrow \\ U_{11} \\ U_{21} \\ \vdots \\ U_{K1} \end{pmatrix}_{K,1} \left(\begin{array}{c} \downarrow \\ X_{11} \quad X_{21} \quad \cdots \quad X_{n1} \end{array} \right)_{1,n}$$

$$U_1 X_1^T = \begin{pmatrix} U_{11}X_{11} & U_{11}X_{21} & U_{11}X_{n1} \\ U_{21}X_{11} & U_{21}X_{21} & U_{21}X_{n1} \\ \vdots & \vdots & \vdots \\ U_{K1}X_{11} & U_{K1}X_{21} & U_{K1}X_{n1} \end{pmatrix}_{K,n} \quad \text{--- (5)}$$

similarly,

$$U_2 X_2^T = \begin{pmatrix} U_{12}x_{12} & U_{12}x_{22} & U_{12}x_{n2} \\ U_{22}x_{12} & U_{22}x_{22} & U_{22}x_{n2} \\ \vdots & \vdots & \vdots \\ U_{K2}x_{12} & U_{K2}x_{22} & U_{K2}x_{n2} \end{pmatrix}_{K,n} \quad \text{--- (6)}$$

and,

$$U_m X_m^T = \begin{pmatrix} U_{1m}x_{1m} & U_{1m}x_{2m} & U_{1m}x_{nm} \\ U_{2m}x_{1m} & U_{2m}x_{2m} & U_{2m}x_{nm} \\ \vdots & \vdots & \vdots \\ U_{Km}x_{1m} & U_{Km}x_{2m} & U_{Km}x_{nm} \end{pmatrix} \quad \text{--- (7)}$$

$$\therefore U X^T = U_1 X_1^T + U_2 X_2^T + \cdots + U_m X_m^T$$

$$\boxed{U X^T = \sum_{i=1}^m U_i X_i^T} \quad \text{Q.E.D.}$$



HW Assignment 7 (Due by 10:30am on Nov 30)

1 Theory (140 points)

1. [Normalized Kernels, 50 + 25 points]

Let $K : X \times X \rightarrow R$ be a kernel function defined over a sample space X .

- (a) Prove that the function below (a *normalized kernel*) is a valid kernel.

$$\frac{K(x, y)}{\sqrt{K(x, x)K(y, y)}}$$

- (b) Why it would not make sense to normalize a Gaussian kernel?

- (c) [Bonus] Which types of input data would benefit from normalizing the kernel function? Explain why, and provide real world examples.

2. [Support Vector Machines, 50 points]

Prove that the sum of slacks $\sum \xi_n$ from the objective function of the SVM formulation with soft margin is an upper bound on the number of misclassified training examples.

3. [Max Margin Hyperplanes, 20 points]

Consider the constrained optimization SVM problem for the separable case shown on slide 12. Show that, if the 1 on the right-hand side of the inequality constraint is replaced by some arbitrary constant $\gamma > 0$, the resulting maximum margin hyperplane is unchanged.

4. [Kernel Techniques, 20 + 20 points]

(a) Show that if $k_1(\mathbf{x}, \mathbf{y})$ and $k_2(\mathbf{x}, \mathbf{y})$ are valid kernel functions, then $k(\mathbf{x}, \mathbf{y}) = k_1(\mathbf{x}, \mathbf{y})k_2(\mathbf{x}, \mathbf{y})$ is also a valid kernel.

(b) (*) Show that if A is a symmetric positive semidefinite matrix, then $k(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T A \mathbf{y}$ is a valid kernel. Hint: Inderjit Dhillon's Linear Algebra Background describes some useful properties of symmetric positive semidefinite matrices.

5. [Positive Definite Matrices (*), 20 points]

Show that a diagonal matrix \mathbf{W} whose elements satisfy $0 < W_{ii} < 1$ is positive definite. Show that the sum of two positive definite matrices is itself positive definite.

6. [Large Margin Perceptron (*), 30 points]

Let \mathbf{u} be a current vector of parameters and \mathbf{x} and \mathbf{y} two training examples such that $\mathbf{u}^T(\mathbf{x} - \mathbf{y}) < 1$. Use the technique of Lagrange multipliers to find a new vector of parameters \mathbf{w} as the solution to the convex optimization problem below:

minimize:

$$J(\mathbf{w}) = \frac{1}{2}\|\mathbf{w} - \mathbf{u}\|^2$$

subject to:

$$\mathbf{w}^T(\mathbf{x} - \mathbf{y}) \geq 1$$

2 Text Classification (100 points)

Train and test the SVM algorithm on the *Spam vs. Non-spam* and *Atheism vs. Religion* classification problems, using the datasets provided for the previous assignment. Use a linear kernel, with the cost parameter $C = 5$. Report and compare the accuracy of the trained SVM models with the perceptron and average perceptron accuracies from the previous assignment.

3 Digit Recognition (200 points)

In this exercise, you are asked to run an experimental evaluation of SVMs and the perceptron algorithm, with and without kernels, on the problem of classifying images representing digits.

1. The UCI Machine Learning Repository at www.ics.uci.edu/~mlearn maintains datasets for a wide variety of machine learning problems. For this assignment, you are supposed to work with the Optical Recognition of Handwritten Digits Data Set. The webpage for this dataset is at:

<http://archive.ics.uci.edu/ml/datasets/Optical+Recognition+of+Handwritten+Digits>

The actual dataset is located at:

<http://archive.ics.uci.edu/ml/machine-learning-databases/optdigits/>

Read the description of the dataset. Download the training set `optdigits.tra` and the test set `optdigits.tes`. Use the first 1000 examples in `optdigits.tra` for *development* and the rest of 2823 examples for *training*. Use all 1797 examples in `optdigits.tes` for *testing*. Scale all the features between [0, 1], as discussed in class, using the min and max computed over the training examples. Create training files for each of the 10 digits, setting the class to 1 for instances of that digit, and to -1 for instances of other digits, i.e. *one-vs-rest* scenario.

2. Train first the linear perceptron, with the number of epochs set to $T \in \{1, 2, \dots, 20\}$. After training each linear perceptron, normalize the learned weight vector. Select for T the value that obtains the best overall accuracy on the development data, and use this value for the remaining perceptron experiments.

Run experiments with the linear and kernel perceptron algorithms. For the kernel perceptron, experiment with polynomial kernels $k(\mathbf{x}, \mathbf{y}) = (1 + \mathbf{x}^T \mathbf{y})^d$ with degrees $d \in \{2, 3, 4, 5, 6\}$, and with Gaussian kernels $k(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|^2 / 2\sigma^2)$ with the width $\sigma \in \{0.1, 0.5, 2, 5, 10\}$. For each hyper-parameter value, you will have trained 10 models, one for each digit. In order to compute the label for a test or development example, you will run the 10 trained models and output the label that obtains the highest score. Compute the accuracy on the development data and identify the hyper-parameter value that obtains the best accuracy. Use the tuned hyper-parameter (d for poly-kernel, σ for Gaussian) to compute the overall performance on the test data.

For each of the three perceptrons (linear, poly kernel and Gaussian kernel) report the total training time, the overall accuracy, and the number of support vectors. Show and compare the corresponding 4 confusion matrices. Which digit seems to be the hardest to classify? Which perceptron / kernel combination achieves the best performance? Which algorithms are slower at training time, and why?

- Run the same experiments using SVMs instead of perceptrons, i.e. linear SVMs and SVMs with polynomial and Gaussian kernels. Use the same tuning scenarios for the hyper-parameters of the polynomial and Gaussian kernels. Use $C = 1$ in all SVM experiments. Report the same types of results and analysis as above, and compare with the perceptron results.

4 Tools

You are free to use MATLAB, R, or packages written in C++/Java/Python such as SVM-LIGHT (C), LIBSVM (C++, Java), or SCIKIT-LEARN (Python) to complete the implementation part of this assignment. Their web sites contain plenty of documentation on how to use them. If you use SCIKIT-LEARN, the following functionality from the `sklearn.svm` will be useful:

- `SVC()`: This is the main class used for SVM classification models. Its implementation is based on LIBSVM. Make sure that you properly map the SVM hyper-parameters to the parameters in the constructor of this class. For example, the `gamma` parameter in the constructor corresponds to our $1/2\sigma^2$ coefficient in the Gaussian kernel. The formulas for the kernels implemented by SVC are described in this User Guide.
- `decision_function(x)`: Once the classifier is trained, this will compute the distance between a sample `x` and the decision hyperplane. This is the quantity that you can use to determine the highest scoring class when training the 10 *one – vs – rest* classifiers: once a classifier is trained for all 10 digits, given a sample `x` you compute this quantity for all 10 classifiers and select the class that corresponds to the classifier with largest decision function value.
- `fit()`: This is the function used to train the classifier.
- `predict(x)`: This is used to calculate the (binary) label for sample `x`.

LIBSVM, and therefore SCIKIT-LEARN too, already implement the *one-vs-rest* classification scheme. In this scheme, you can directly use the training dataset with the 10 original labels, and SCIKIT-LEARN will train the 10 binary classifiers for you. You can use this capability for this assignment, however bonus points will be given if you train the 10 binary classifiers directly, as described for the perceptron algorithm above, by creating a binary training dataset for each class.

5 Submission

Turn in a hard copy of your homework report at the beginning of class on the due date. Electronically submit on Blackboard a `hw07.zip` file that contains the `hw07` folder in which you place the code and the datasets. Make sure you include a `README.txt` file explaining how the code is supposed to be used to replicate the results included in the report. The screen output produced when running the code should be redirected to (saved into) an `output.txt` file.

On a Linux system, creating the archive can be done using the command:

```
> zip -r hw07.zip hw07
```

Please observe the following when handing in homework:

1. Structure, indent, and format your code well.
2. Use adequate comments, both block and in-line to document your code.
3. **Do not submit third-party ML packages on Blackboard!** Just explain in the REAMDE file how you use external packages.
4. Make sure your code runs correctly when used in the directory structure shown above.
5. **Type and nicely format the project report**, including discussion points, tables, graphs etc. so that it is presentable and easy to read.
6. Working code and/or correct answers is only one part of the assignment. The project report, including discussion of the specific issues which the assignment asks about, is also a very important part of the assignment. Take the time and space to make an adequate and clear project report. On the non-programming learning-theory assignment, clear and complete explanations and proofs of your results are as important as getting the right answer.

HW 07

Bhishan Poudel

63/50

Lecture

Normalized Kernel is a valid kernel.

Defn: Let $K(x,y)$ is a valid kernel, then, there exist a feature map $\phi(\cdot)$ in some Hilbert space that,

$$K(x,y) = \phi(x) \cdot \phi(y)$$

$$\text{Now, } \hat{K}(x,y) = \frac{K(x,y)}{\sqrt{K(x,x), K(y,y)}}$$

$$= \frac{\phi(x) \cdot \phi(y)}{\|\phi(x)\| \|\phi(y)\|} = \hat{\phi}(x) \cdot \hat{\phi}(y)$$

where $\hat{\phi}(x) = \frac{\phi(x)}{\|\phi(x)\|}$

since $K(x,y)$ is kernel $\Rightarrow \phi(x) \cdot \phi(y) \geq 0$

also, denominator $\Rightarrow \|\phi(x)\| \|\phi(y)\| \geq 0$

$\therefore \hat{K}(x,y)$ is non-negative. } $\therefore \hat{K}(x,y)$
 also $\hat{K}(x,y)$ is symmetric. } is a
 valid kernel.

Q17b Gaussian function,

Gaussian fn

$$g(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{\|x-\mu\|^2}{2\sigma^2}}$$

(normalized gaussian function when $y=\mu$)

The gaussian kernel is,

Gaussian Kernel

$$K(x,y) = e^{-\frac{\|x-y\|^2}{2\sigma^2}} = e^{-\gamma \|x-y\|^2}$$

where $\gamma = \frac{1}{2\sigma^2}$

$$K(0,0) = e^0 = 1$$

$K(0,0) = 1$ for gaussian kernel.

Since gaussian kernel is already normalized ($K(0,0)=1$), it would make no sense to normalize gaussian kernel.

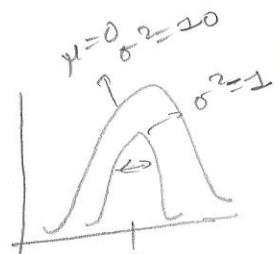
3C some of the kernels used in practice

linear kernel $K(x,y) = \alpha^T y + c = \langle \alpha, y \rangle + c = \vec{\alpha} \cdot \vec{y} + c$

$\alpha = \text{slope}$
 $d = \text{degree}$
 $c = \text{const}$

polynomial kernel $K(x,y) = (\alpha^T y + c)^d$

Gaussian Kernel $K(x,y) = e^{-\frac{1}{2\sigma^2} \|x-y\|^2} = e^{-\pi/\sigma^2 \|x-y\|^2}$



$$\text{FWHM} = 2\sqrt{2\ln 2} \cdot \sigma$$

(σ or γ determines the width of the gaussian curve.
 overestimate \rightarrow behave like linear kernel
 underestimate \rightarrow lack regularization and high sensitive to noise
 tuning of $\sigma \rightarrow$ very important, otherwise very low performance

exponential kernel $K(x,y) = e^{-\frac{1}{\sigma} \|x-y\|}$

sigmoid kernel $K(x,y) = \tanh(\alpha^T y + c)$

Gaussian kernels are already normalized and we need to normalize polynomial kernels.

cosine normalization

- Reduces dimension of data by 1 (since it projects all the data inside a unit radius sphere)
- Good for high dimensional input data
e.g. optdigits with 64 features
- NOT good for low dimensional data
e.g. Iris data with only two four features.

(ON2)
GO
GO

prove! sum of slack $\geq \epsilon_n$ is an upper bound on the number of misclassified training examples.

solution:

Input data

X

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_i \\ \vdots \\ x_N \end{bmatrix}$$

Target

t

$$\begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_i \\ \vdots \\ t_N \end{bmatrix}$$

$$t \in \{-1, 1\}$$

$$x_i \in \mathbb{R}^d$$

$$x_i \in \{x_{i1}, x_{i2}, \dots, x_{id}\} \quad d = \text{dimension of features.}$$

$$\text{hypothesis} \quad h_i = w^T x_i + b$$

$$t_i \in \{-1, 1\}$$

algorithm :

$$\xrightarrow{\text{constraint}} t_i (w^T x_i + b) \geq 1 - \epsilon_i$$

$\xrightarrow{\text{minimization}}$
objective

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

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

$$\text{with constraint } t_i (w^T x_i + b) \geq 1 - \epsilon_i$$

consider

$$t_i = +1$$

(correct classification)

$$w^T x_i + b > 0$$

$$t_i (w^T x_i + b) > 0 \quad \because t_i = +1$$

~~constraint~~ $t_i (w^T x_i + b) \geq 1 - \xi_i \quad ; \quad \xi_i \geq 0$

aside:

for perceptron

$$h_i = \text{sgn}(w^T x_i + b)$$

$$h_i = t_i (w^T x_i + b)$$

and if

if $t_i h_i$ we update w and b

~~$t_i (w^T x_i + b) \geq 1 - \xi_i$~~

$$+ve \geq 1 - \xi_i$$

since constraint want $\xi_i \geq 0$

~~ξ_i cannot be greater than 1~~
 $0 \leq \xi_i \leq 1$

$$t_i (w^T x_i + b) \geq 1 \text{ if } \xi_i = 0$$

(hard margin)

consider $t_i = +1$ (correct classification)

$$w^T x_i + b < 0$$

$$t_i (w^T x_i + b) < 0 \quad \because t_i = +1$$

constraint:

$$t_i (w^T x_i + b) \geq 1 - \xi_i$$

$$+ve \geq 1 - \xi_i$$

$$\xi_i \geq 1 - (+ve)$$

$$0 \leq \xi_i \leq 1$$

$$\text{and } t_i (w^T x_i + b) \geq 1 \text{ if } \xi_i = 0$$

\Rightarrow for correct classification

$$0 \leq \xi_i \leq 1$$

consider $t_i = +1$

(misclassification)

$$w^T x_i + b < 0$$

$$t_i (w^T x_i + b) < 0 \quad \because t_i = +1$$

constraint:

$$t_i (w^T x_i + b) \geq 1 - \xi_i$$

$$-ve \geq 1 - \xi_i$$

$$\xi_i \geq 1 + (-ve)$$

$$\xi_i > 1$$

$$t_i = -1 \quad (\text{misclassification})$$

$$w^T x_i + b > 0$$

$$t_i (w^T x_i + b) < 0 \quad \because t_i = -1$$

constraint:

$$t_i (w^T x_i + b) \geq 1 - \xi_i$$

$$-ve \geq 1 - \xi_i$$

$$\xi_i \geq 1 + (+ve)$$

$$\xi_i > 1$$

\Rightarrow for misclassification

$$\xi_i > 1$$

Now, for correct classification

$$e_{ii} \geq 0$$

$$\sum e_{ii} \geq 0$$

i: correct

for misclassified

$$e_{ij} > 1$$

$$\sum e_{ij} > \sum 1$$

i: misclassified j: misclassified

$$\sum e_{ii} > \# \text{ of misclassified}$$

$$\therefore \sum_{i \in \text{correct}} e_{ii} + \sum_{i \in \text{misclassified}} e_{ii} > \# \text{ of misclassified}$$

$$\therefore \boxed{\sum_{i=1}^N e_{ii} > \# \text{ of misclassified}}$$

Q.E.D.

This shows that sum of slack (e_{ii}) over all the N training examples is greater than total number of misclassified examples, hence, sum of slacks is an upper bound to the number of misclassified examples.

Aside: the constrained optimization of SVM algorithms

$t_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$ is called Hard margin SVM, and,

$t_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - e_{ii}$ is called Soft margin SVM.

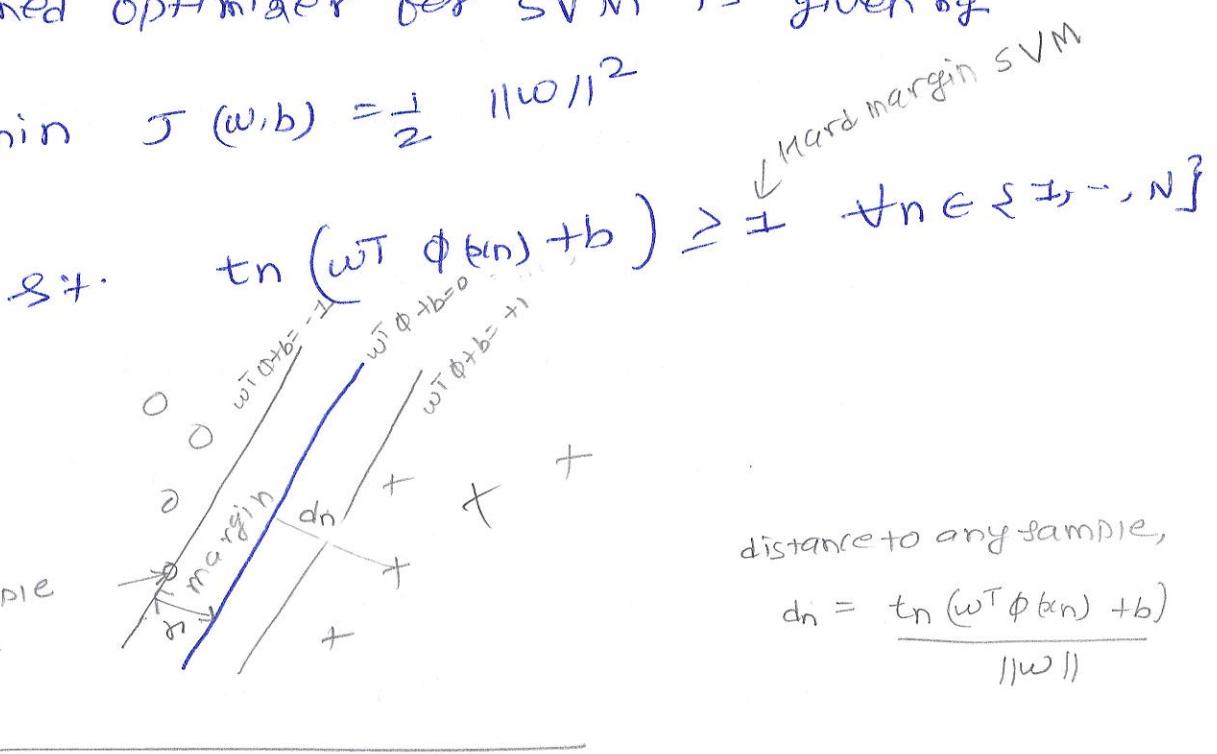
They both are non-trivial algorithm (unlike vanilla perceptron) and use Lagrange multipliers to solve the optimization.

QN3

max margin hyperplane

for the linearly separable dataset, the constrained optimizer for SVM is given by

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



from slides,
(lecture 07, page 5)

$$\text{margin} = \min_n \frac{t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b)}{\|\mathbf{w}\|}$$

we find parameters w and b that maximizes margin,

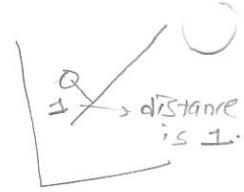
$$w^*, b^* = \underset{w, b}{\operatorname{argmax}} \cdot \min_n t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) / \|w\|$$

Since w^*, b^* are obtained from argmax, they do not depend on rescaling of w, b .

So, without loss of generality, we can say that the nearest distance to the sample from hyperplane is 1.

i.e. for closest point(s),

$$t_n(\omega^T \phi(x_n) + b) = 1$$



then, our optimizing constraint becomes,

$$t_n(\omega^T \phi(x_n) + b) \geq 1 \quad \forall n \in \{1, \dots, N\}$$

[Aside: previously we used

$$t_n(\omega^T \phi + b) = 1$$

$$t_n h_n = 1$$

(for perceptron)

$$\begin{cases} \omega^T \phi(x_n) + b > 0 & \text{if this is +} \\ \omega^T \phi(x_n) + b < 0 & \text{if this is -} \\ h_n = \text{sgn}(\omega^T \phi(x_n) + b) \end{cases}$$

now, with inclusion of this constraint, the new optimization problem for SVM becomes,

$$\left\{ \begin{array}{l} \text{minimize} \\ \sigma(\omega | b) = \frac{1}{2} \|\omega\|^2 \end{array} \right.$$

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

here we have
use the margin
 $\gamma = 1$.

we have to show that if $\gamma_1 = 0$, then
also the decision hyperplane is unchanged.

Solution :

The distance to the n th sample
from the decision hyperplane is,

$$d_n = \frac{t_n(\omega^T \phi(t_n) + b)}{\|\omega\|} \quad \} \rightarrow ①$$

If we rescale the weight vector to ± 1 ,
 $\|\omega\| = 1$, then,

$$d_n = t_n(\omega^T \phi(t_n) + b) \quad \} \rightarrow ②$$

then, maximizing margin algorithm becomes,

$$\max_{\omega, b} \min_n d_n \quad \} \rightarrow ③$$

s.t. $\|\omega\| = 1$

Let γ_1 be the maximum margin to ^{closest} training example,

$$\forall n \quad t_n(\omega^T \phi(t_n) + b) \geq \gamma_1$$

$$\Rightarrow \quad t_n\left(\frac{\omega^T \phi}{n} + \frac{b}{n}\right) \geq 1$$

$$\Rightarrow \quad t_n(N^T \phi + B) \geq 1 \quad \} \rightarrow ④$$

where $v = \frac{w}{\|w\|}$, $B = \frac{b}{\|w\|}$, γ = margin to the closest example.

then, $v = \frac{w}{\|w\|}$

$$\|v\|^2 = \left\| \frac{w}{\|w\|} \right\|^2 = \frac{1}{\|w\|^2} \|w\|^2$$

$$\boxed{\|v\|^2 = \frac{1}{\|w\|^2}} \quad \text{since } \|w\| = 1$$

Hence, our original optimization problem,

$$\max_{w, b} \min_n d_n = \min_n (w^T \phi(x_n) + b)$$

reduces to maximizing $\frac{1}{\|w\|^2}$

i.e. minimizing $\|v\|^2$

i.e. minimizing $\frac{1}{2} \|v\|^2$

which we can write as,

$$\min \frac{1}{2} \|v\|^2$$

$$\text{s.t. } \min_n (v^T \phi(x_n) + b) \geq 1.$$

(Eqn ④)

which is same as

$$\min \mathcal{J}(w, b) = \frac{1}{2} \|w\|^2$$

$$\text{s.t. } \min_n (w^T \phi(x_n) + b) \geq 1 \quad \text{in Eqn ④}$$

A \rightarrow D

~~given~~ \Rightarrow product of two kernels is a kernel.

Let ϕ_1, ϕ_2 are the feature map of kernels k_1, k_2

then,

$$\begin{aligned} K_1(x_1, x_2) \cdot K_2(b_1, b_2) &= (\phi_1(x_1) \quad \phi_1(x_2)) \quad (\phi_2(b_1) \quad \phi_2(b_2)) \\ &= \left(\sum_{i=1}^{\infty} f_i(x_1) f_i(b_2) \right) \left(\sum_{j=1}^{\infty} g_j(b_1) g_j(x_2) \right) \\ &= \sum_{i=1}^{\infty} \sum_{j=1}^{\infty} f_i(x_1) f_i(b_2) g_j(b_1) g_j(x_2) \\ &= \sum_{i=1}^{\infty} \sum_{j=1}^{\infty} (f_i(x_1) g_j(x_1)) \quad (f_i(b_2) g_j(b_2)) \\ &= \sum_{i=1}^{\infty} \sum_{j=1}^{\infty} h_{ij}(x_1) h_{ij}(b_2) \end{aligned}$$

where $h_{ij}(b)$ is a feature vector in feature space ϕ_3 such that

$$h_{ij}(x) = f_i(x) g_j(x)$$

$$\therefore K_1(x_1, x_2) \cdot K_2(b_1, b_2) = \phi_3(x_1)^T \phi_3(x_2)$$

is a valid kernel.

~~X~~ ~~any b~~ $x^T A y$ is valid kernel if A is sym PSD matrix.

join: Given $A \rightarrow$ symmetric positive semidefinite

$A^T A$ is also sym PSD (from
inherit Shallow's reference)

Now,

$$\text{Let, } K(x,y) = x^T A^T A y$$

$$= (Ax)^T (Ay) \quad \because (AB)^T = B^T A^T \\ = \phi(x)^T \phi(y)$$

where we define feature map
 $\phi(z) = Az$

$\therefore K(x,y)$ is a valid kernel.

~~ANSWER~~ Diagonal matrix with positive elements
is positive definite matrix.

~~20/20~~
SOLN: Let $w = w_{n,n} \in \mathbb{R}^{n \times n}$ is a diagonal matrix then,

$$w = \begin{bmatrix} w_{11} & & & \\ & w_{22} & & \\ & & \ddots & \\ & & & w_{nn} \end{bmatrix}_{n,n}$$

Let \vec{x} be a n -dimensional column vector,

$$\vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}_{n,1}$$

then, the inner product of diagonal matrix w associated with \vec{x} is,

$$\vec{x}^T w \vec{x} = [x_1 \ x_2 \ \dots \ x_n]_{1,n} \begin{bmatrix} w_{11} & & \\ & \ddots & \\ & & w_{nn} \end{bmatrix}_{n,n} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}_{n,1}$$

$$= \sum_{j=1}^n w_j x_j^2$$

where $w_j = w_{jj}$
(diagonal elements)

$$x^T W x = \sum_{j=1}^n w_j x_j^2 = \text{positive number}$$

here all the diagonal elements w_j are positive
 $(0 < w_{ii} \leq +\infty)$.
 then the RHS quantity is positive
 and hence W is a positive definite matrix.

After: $W = \text{diag}(w_{11}, w_{22}, \dots, w_{nn})$, $w_{ii} \geq 0$
 let $\sqrt{W} = \text{diag}(\sqrt{w_{11}}, \sqrt{w_{22}}, \dots, \sqrt{w_{nn}})$

$$\begin{aligned} x^T W x &= (x^T \sqrt{W}^T)(\sqrt{W} x) \\ &= (\sqrt{W} x)^T (\sqrt{W} x) \\ &= \|\sqrt{W} x\|^2 \geq 0 \end{aligned}$$

$\therefore W$ is PSD.

~~ANSWER~~ sum of two PSD matrices is PSD matrix.

Let W is a positive ~~semi-definite~~ matrix
 $\Rightarrow \alpha^T W \alpha \geq 0 \quad \forall \alpha \in \mathbb{R}^n$

Let V is a positive semi-def matrix
 $\Rightarrow \alpha^T V \alpha \geq 0$

now, $\alpha^T W \alpha + \alpha^T V \alpha \geq 0$

or, $\alpha^T (W+V) \alpha \geq 0$ (distributive property of matrix)

this means that sum of two matrices $W+V$ is also positive semi-definite.

proved

Ques $\omega = u + \lambda(b-y)$

\rightarrow s.t. $J(\omega) = \frac{1}{2} \|\omega - u\|^2$

$$= \frac{1}{2} \|\lambda(b-y)\|^2$$

constraint: $\omega^T(b-y) = u^T(b-y) + \lambda \|b-y\|^2 \geq 1$

for, $\lambda \|b-y\|^2 \geq 1 - u^T(b-y)$

$$\lambda \geq \frac{1 - u^T(b-y)}{\|b-y\|^2}$$

Third, optimization problem:

$$\min_{\lambda} J(\lambda) = \frac{1}{2} \lambda^2 \|b-y\|^2$$

s.t. $\lambda \geq \frac{1 - u^T(b-y)}{\|b-y\|^2}$

Again, $J(\lambda)$ will be minimum when λ will be minimum.
 But $\min \lambda = \frac{1 - u^T(b-y)}{\|b-y\|^2}$

so, $\lambda = \frac{1 - u^T(b-y)}{\|b-y\|^2}$

then, the ω is given by
 (eliminate λ)

$$\boxed{\omega = u + \frac{1 - u^T(b-y)}{\|b-y\|^2} (b-y)}$$

Ans

Lag6

Lagrange multipliers

Here, the optimization problem is,

minimize w $J(w) = \frac{1}{2} \|w - u\|^2$

s.t. $w^T (x-y) \geq 1$

$$w^T x - w^T y \geq 1$$

$$1 - w^T x + w^T y \leq 0$$

primal Lagrangian

$$L_p(w, \alpha, y) = \frac{1}{2} (w-u)^2 + \alpha (1 - w^T x + w^T y)$$

$\alpha \geq 0$ is Lagrange multiplier

Set gradients to zero

$$\frac{\partial L_p}{\partial w} = w - u - \alpha x + \alpha y \Rightarrow w = u + \alpha x - \alpha y = u + \alpha(x-y)$$

$$\frac{\partial L_p}{\partial \alpha} = 1 - w^T x + w^T y \Rightarrow 1 = w(x-y) \Rightarrow \alpha = \frac{1}{w}$$

$$\text{then, } w = u + \frac{\alpha}{\alpha} = u + \frac{1}{w} \Rightarrow w^2 = uw + 1 \Rightarrow w^2 - uw - 1 = 0$$

$$w^* = \frac{u \pm \sqrt{u^2 + 4}}{2}$$

This is
not
solution
look
this

In gradient descent algorithm we initialize $\vec{u} = 0$

and $\alpha > 0$ and find the optimum value of w^* .

Bhishan Poudel

HW Assignment 8 (Due by 10:30am on Dec 7)

1 Theory (150 points)

1. [Kernel Nearest Neighbor, 50 points]

The nearest-neighbour classifier 1-NN assigns a new input vector \mathbf{x} to the same class as that of the nearest input vector \mathbf{x}_n from the training set, where in the simplest case, the distance is defined by the Euclidean metric $\|\mathbf{x} - \mathbf{x}_n\|^2$. By expressing this rule in terms of scalar products and then making use of kernel substitution, formulate the nearest-neighbour classifier for a general nonlinear kernel.

2. [Distance-Weighted Nearest Neighbor, 50 points]

We have seen how to use kernels to formulate a distance-weighted nearest neighbor algorithm, when the labels are binary. Formulate a kernel-based, distance-weighted nearest neighbor that works for K classes, where $K \geq 2$.

3. [Naive Bayes, 50 points]

The Naive Bayes algorithm for text categorization presented in class treats all sections of a document equally, ignoring the fact that words in the title are often more important than words in the text in determining the document category. Describe how you would modify the Naive Bayes algorithm for text categorization to reflect the constraint that words in the title are K times more important than the other words in the document for deciding the category, where K is an input parameter (include pseudocode).

4. [Logistic Regression (*), 50 points]

Assume that a binary feature x_i is equal to 1 for all training examples \mathbf{x} belonging to a particular class C_k , and zero otherwise (i.e. x_i perfectly separates examples from class C_k from all other examples). Show that in this case the magnitude of the ML solution for \mathbf{w}_k goes to infinity, thus motivating the use of a prior over the parameters (Hint: use the fact that the gradient on slide 24 must vanish at the solution).

2 Submission

Turn in a hard copy of your homework report at the beginning of class on the due date. On this theory assignment, **clear and complete explanations and proofs of your results are as important as getting the right answer.**

HW08

Bishan Poudel

K-Nearest Neighbor (Kernel based classifier)

① 1 nearest neighbor

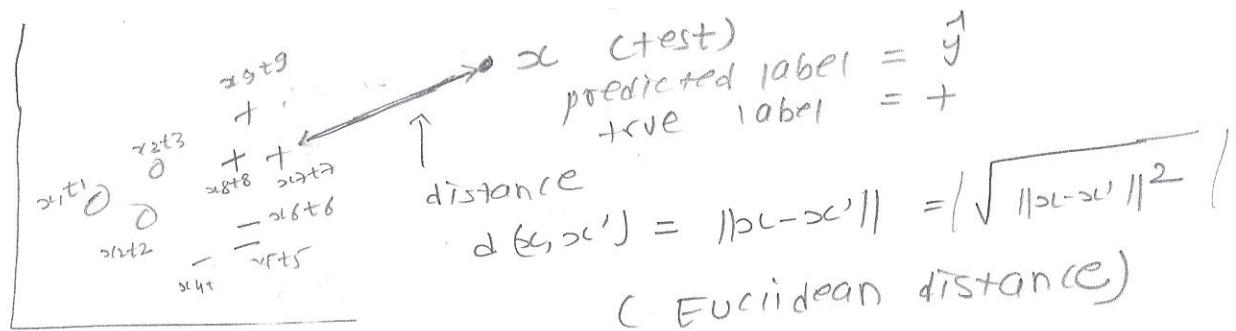
training dataset : $(x_1, t_1), (x_2, t_2), \dots, (x_n, t_n)$

test dataset : x

required output: \hat{y} (label of x or class of x)

t_1, t_2, \dots, t_n are classes from c_1, c_2, \dots, c_m

for example {circle plus minus }
 $c_1 \quad c_2 \quad c_m$



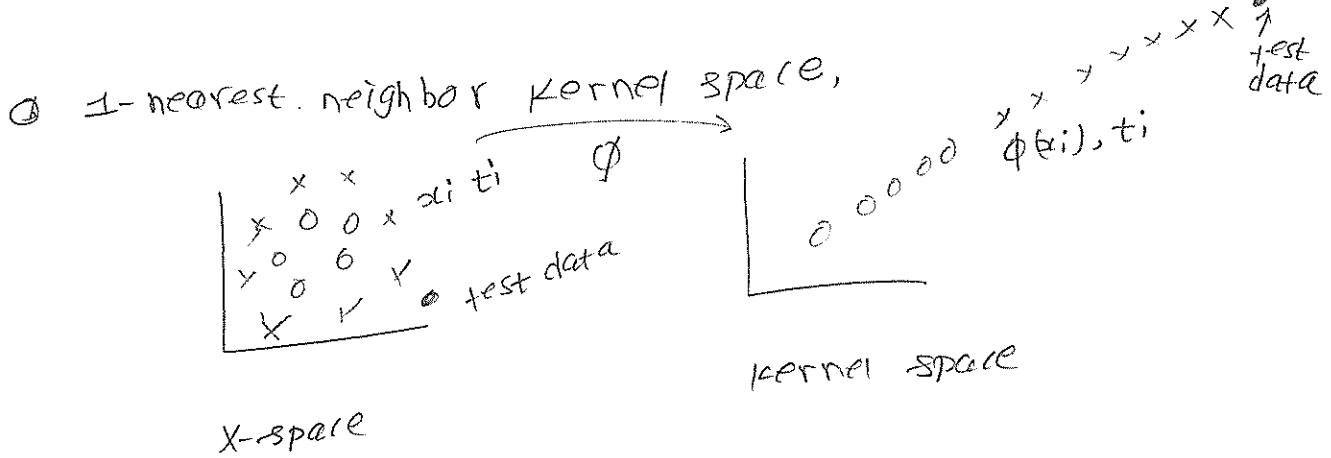
How to find \hat{y} ?

⇒ For 1-nearest neighbor, calculate the distance (maybe Euclidean, cosine, edit, etc) from test point \vec{x} to all the data points $x_i, 1 \leq i \leq n$ then find the minimum distance data point x' .

Assign label of x' as the predicted label to \vec{x} .

$$\hat{y} = \text{class of } x'$$

= class of x_i such that $d(x_i)$ is minimum



Here, when calculating distance, we use kernel method.

X-space: $d(x, x') = \|x - x'\| = \sqrt{\|x - x'\|^2} = \sqrt{(x - x')^T (x - x')}$

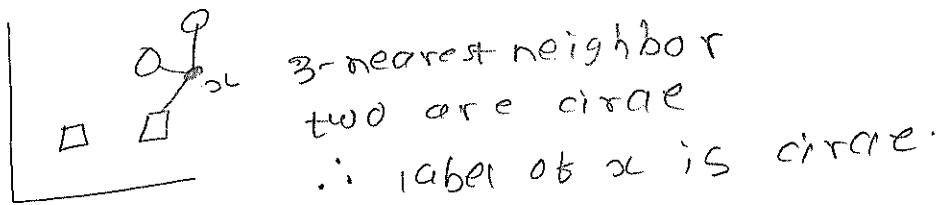
(Euclidean distance)

Kernel-Space: $d(\phi(x), \phi(x')) = \sqrt{(\phi(x) - \phi(x'))^T (\phi(x) - \phi(x'))}$
 $= \sqrt{K(x, x')}$

→ we compute the distance of test \vec{x} to the all the training examples x_i ; using kernel method then choose the minimum distance training point x' and then,

$$\text{label of } x = \text{label of } x'$$

⑥ K-nearest neighbor kernel method.



Kernel \Rightarrow distances are calculated using
kernel trick

$$d(x_i, x_j) = \sqrt{k_{\phi}(x_i, x_j)}$$

Kernel can be polynomial kernel,
gaussian kernel and so on.

Let x_1, x_2, \dots, x_K are nearest of x

then,

$$y(x) = \underset{t \in T}{\operatorname{argmax}} \sum_{i=1}^K \delta_t(t_i)$$

Hence, $t_i = \begin{cases} 0, & \text{circle square} \\ 1, & \text{square} \end{cases}$

$K=3$

for x_1 and x_2 $t=t_i=0$ number = 2

for x_3 $t=t_i=1$ number = 1

$$\operatorname{argmax} = 0$$

Summary

input $(x_1, t_1), (x_2, t_2), \dots, (x_n, t_n)$ $t_n \in \{t_1, t_2, \dots, t_k\}$

test \sim , $t = ?$

map input from x -space to kernel-space

x -space input $\phi(x_1, t_1) \quad \phi(x_2, t_2) \quad \dots \quad \phi(x_n, t_n)$

kernel space test $\phi(x), t = ?$

compute similarity $K(\phi(x), \cdot)$ for all ϕ_i , align

(compute distance $= \|\sqrt{K(\phi(x), \cdot)}\|$ is redundant since
distance \leq similarity)

then predicted class of x ,

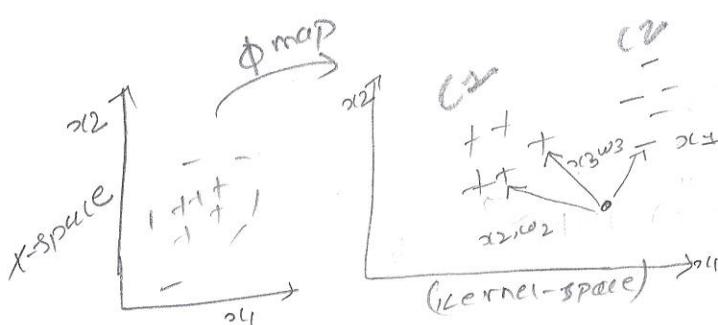
$$y(x) = \arg \max_{t \in T} \sum_{i=1}^k \delta_t(t_i)$$

[PN2] DISTANCE-weighted Nearest Neighbor for classification

@ Kernel space

@ distance-weighted KNN binary classification

ϕ map



$$t_i \in \{-1, +1\}$$

3-nearest neighbor

contribution from + : $w_2 + w_3$

contribution from - : w_1

suppose $w_2 + w_3 > w_1$

then label of $x = +$

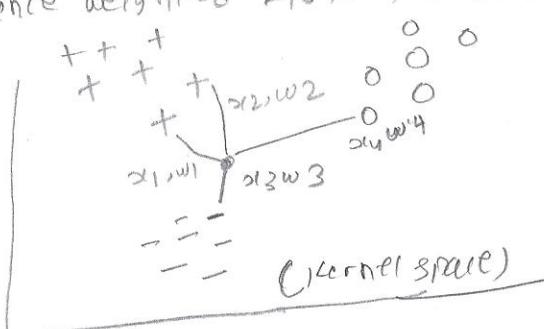
- x' is test example
- x_1, x_2, \dots, x_k are the nearest points to x' , calculated using Euclidean distance in kernel space.
- $w_k(x_i)$ is the measure of similarity between point x' and x_i in kernel space.

$$\text{Weighted sum} = \sum_{i=1}^k w_k(x_i) \delta_t(t_i)$$

$$y = \text{sign} \left(\sum_{i=1}^k w_k(x_i) \delta_t(t_i) \right)$$

- (b) kernel-based distance weighted KNN for multiclass classification

Note:
nearest points
have more
contribution



4-nearest neighbor

contribution from class $c_1 = +$
 $= w_1 + w_2$

contribution from class $c_2 = -$
 $= w_3$

$$y = \operatorname{argmax}_{t \in T} \sum_{i=1}^k w_k(x_i) \delta_t(t_i)$$

contribution from class $c_3 = 0$
 $= w_4$

suppose $w_1 + w_2 > w_3 \quad \therefore y = +$
and $w_1 + w_2 > w_4 \Rightarrow$

Ans

Q3 Naïve Bias

① simple text classification using naïve bias

→ each document is an example vector

→ a document x has n words w_1, w_2, \dots, w_n

→ from all the documents $D = \{x_1, x_2, \dots, x_m\}$

we create vocabulary file (this is feature vector for whole dataset D)

→ vocabulary $V = \{w_1, w_2, \dots, w_{|V|}\}$, there are $|V|$ number of words in vocabulary

→ each document x belongs to one of the categories

$C_K = \{c_1, c_2, \dots, c_K\}$

for example $D = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} \text{line one} \\ \text{line two} \\ \text{line three} \end{bmatrix}$

input data target (category label)

news
sports
religion

we want \Rightarrow given a document x find its category
 $p(c_k|x)$ probability of document x belonging to category c_k (e.g. news)

Naïve Bias gives \Rightarrow $\begin{cases} \textcircled{1} p(c_k) & \text{probability of each categories (priors)} \\ \textcircled{2} p(w_i|c_k) & \text{conditional probability of data } x \text{ w.r.t. category } c_k \end{cases}$

then we use Bayes theorem

to get $p(c^*|x)$.

$$p(c^*|x) = \frac{p(x|c^*) p(c^*)}{p(x)}$$

$$= p(c^*) \prod_{i=1}^n p(w_i|c^*)$$

where c^* corresponds to the $p(c^*|x)$ that has maximum probability

$$c^* = \arg \max_{c_k} p(c_k|x)$$

NB for toy classification when all words have equal importance

D: documents x_1, x_2, \dots, x_n

C $\{1, 2, \dots, K\}$: classes

$\{w_1, w_2, \dots, w_N\}$

below threshold words

V: vocabulary (unique words in D) \rightarrow top words

DK: subset examples having class C K

n $_K$: # of words in DK

1. for each category C K :

2. $D_K = \text{subset with category } C_K$

$$P(C_K) = \frac{|D_K|}{|D|} \quad (\text{prior})$$

3.

4.

$n_K = \# \text{ words in } D_K$

for each word $w_i \in V$:

let $n_{wi} = \# \text{ of } w_i \text{ in } D_K$

6.

7.

$$\text{set } P(w_i | C_K) = \frac{n_{wi} + 1}{n_K + N} \quad \begin{array}{l} \text{Laplace} \\ \text{smoothing} \end{array}$$

8. return priors $P(C_K)$
and cond prob $p(w_i | C_K)$

Then we find the posterior probability for each class C K ,

$$P(C_K|x) = \frac{P(x|C_K) \cdot P(C_K)}{P(x)}$$

$$\propto p(C_K) P(x)$$

$$\propto p(C_K) \cdot \prod_{i=1}^N p(w_i | C_K) \quad \begin{array}{l} \text{wi are words in} \\ \text{test document x} \end{array}$$

$$\propto \ln [p(C_K) \cdot \prod_i p(w_i | C_K)]$$

$$\propto \ln p(C_K) + \sum_i \ln (p(w_i | C_K))$$

(to prevent underflow)

Testing

$$c^* = \arg \max_{C_K} p(C_K|x)$$

$$c^* = \arg \max_{C_K} p(C_K) \prod_i p(w_i | C_K)$$

$$\propto \arg \max_{C_K} [\ln p(C_K) + \sum_i \ln (p(w_i | C_K))]$$

Example of Naive Bayes

train	doc	words	class	Total
1	1	chinese Beijing chinese	$c_1 = c$	$n_1 = 8 \text{ words}$ $d_1 = 3 \text{ examples}$ $n_1 + n_1 = 8 + 6 = 14$
	2	chinese chinese shanghai	c_1	
	3	chinese macao	c_1	
	4	TOKYO JAPAN chinese	c_2	
test	1	chinese Chinese Chinese TOKYO Japan (\vec{x})	?	$n_2 = 3 \text{ words}$ $d_2 = 1 \text{ example}$ $n_2 + n_1 = 3 + 6 = 9$ $\exists = \text{total words}$ 6 unique words $ V = 6$

vocabulary: $\{w_1^{\text{chinese}}, w_2^{\text{Beijing}}, w_3^{\text{shanghai}}, w_4^{\text{macao}}, w_5^{\text{TOKYO}}, w_6^{\text{JAPAN}}\} \quad |V| = 6$

word frequency:

word freq in c_1 :

Laplace smoothing:

$\frac{n_{ik} + 1}{n_k + |V|}$:

category $c_1 = C$

prior $P(C) = |D_1| / |D| = 3/4$ (3 documents belongs to class C_1)

number of words in class, $n_1 = 8$

conditional probabilities ($P(w_i | C) = \frac{n_{ik} + 1}{n_k + |V|}$)

chinese $P(w_1 | C) = \frac{6+1}{8+6} = \frac{7}{14} = \frac{1}{2}$ in $P(w_1 | C) =$

Beijing $P(w_2 | C) = \frac{2+1}{14} = \frac{3}{14}$ in $P(w_2 | C) =$

shanghai $P(w_3 | C) = \frac{2}{7}$ in $P(w_3 | C) =$

macao $P(w_4 | C) = \frac{2}{7}$ in $P(w_4 | C) =$

TOKYO $P(w_5 | C) = \frac{0+1}{8+6} = \frac{1}{14}$ in $P(w_5 | C) =$

JAPAN $P(w_6 | C) = \frac{0+1}{14} = \frac{1}{14}$ in $P(w_6 | C) =$

category $c_2 = J$

prior $P(C_2) = P_2 / |D| = 1/4$
(1 document belongs to class C_2)

$n_2 = 3$

conditional probabilities

$$P(w_1 | C_2) = \frac{4+1}{3+6} = \frac{5}{9}$$

$$P(w_2 | C_2) = \frac{0+1}{3+6} = \frac{1}{9}$$

$$P(w_3 | C_2) = \frac{1}{3}$$

$$P(w_4 | C_2) = \frac{1}{3}$$

$$P(w_5 | C_2) = \frac{2}{3}$$

$$P(w_6 | C_2) = \frac{2}{3}$$

③ choosing class

$$P(C_1 | \vec{x}) \propto P(C_1) \prod_{i \in C_1} P(w_i | C_1) = \frac{3}{4} \cdot \left(\frac{2}{3}\right)^3 \cdot \frac{1}{14} \cdot \frac{1}{14} = 0.0003$$

$$P(C_2 | \vec{x}) \propto P(C_2) \prod_{i \in C_2} P(w_i | C_2) = \frac{1}{4} \cdot \left(\frac{1}{3}\right)^3 \cdot \frac{2}{3} \cdot \frac{2}{3} = 0.0001$$

$$C^* = \operatorname{argmax}_{C_k} P(C_k) \prod_{j=1}^n P(w_j | C_k) = \operatorname{argmax}_{C_k} \{0.0003, 0.0001\} = \text{Label 0 or } 0.0003 = C_1$$

P.T.O.

⑥ give more importance to title of doc

① for each category C_k :

find class prior $P(C_k) = |D_{C_k}| / |D|$ (where D_{C_k} is subset of D with class C_k)

② # loop through title words

③ let n_k = number of title words in D_{C_k} all examples

for each title word $w_i \in V_1$:

let n_{ki} = number of words in title of all D_{C_k} examples

find conditional prob of each title words

$$P_1(w_i | C_k) = \frac{n_{ki} + 1}{n_k + N_1} \times K \leftarrow \begin{array}{l} \text{multiplication factor} \\ \text{for title words} \end{array}$$

(4)

loop through body words

let n_k = number of body words

for each body word $w_i \in V_2$:

let n_{ki} = number of words in body of all D_{C_k} examples (for example C_k is class 1)

find conditional prob of each body words

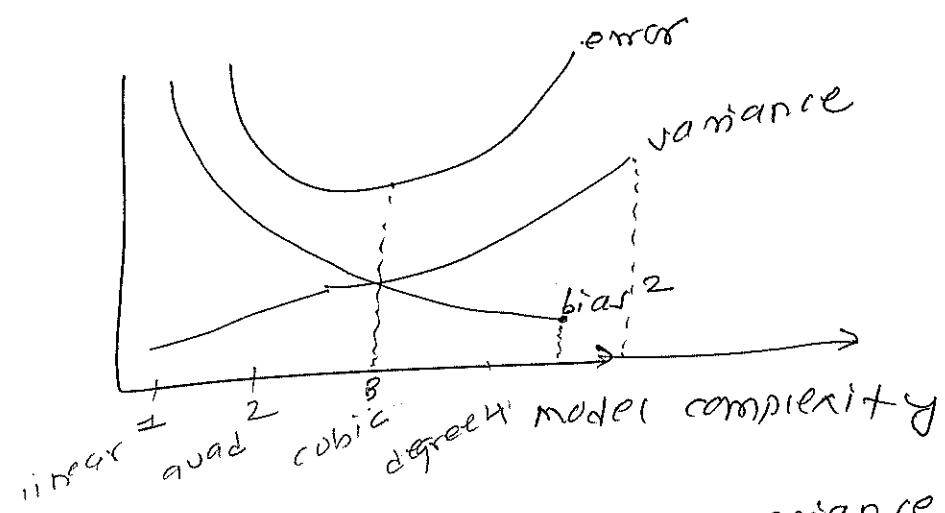
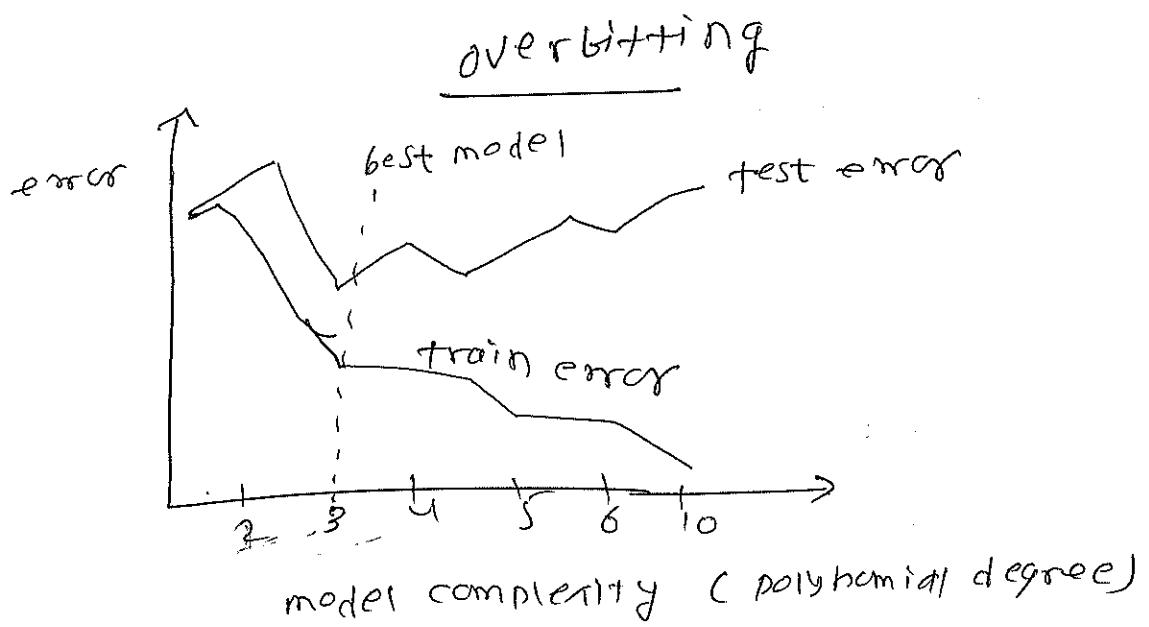
$$P_2(w_i | C_k) = \frac{n_{ki} + 1}{n_k + |V_2|}$$

\Rightarrow output $P(C_k)$, $P_1(w_i | C_k)$, $P_2(w_i | C_k)$

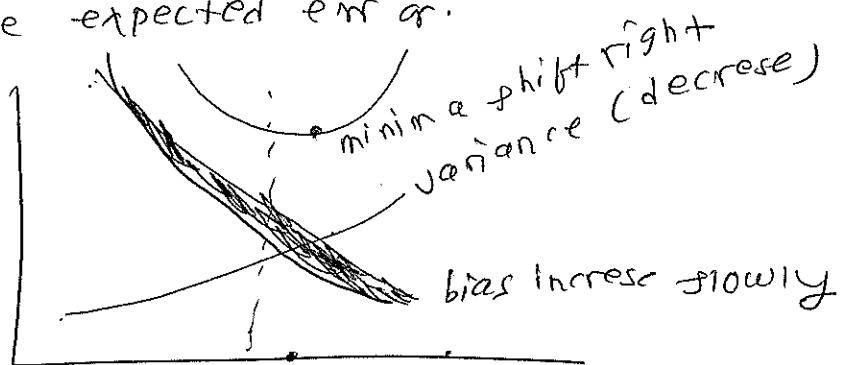
choosing class:

we again suppose title words and body words are independent,

$$P(C_k | \alpha) \propto \underbrace{P_1(\alpha | C_k) P(C_k)}_{\text{for title}} \cdot \underbrace{P_2(\alpha | C_k) P(C_k)}_{\text{for body}}$$



when training data increase, variance is reduced
and slightly more complicated model minimizes
the expected error.



① show that $x^T A y$ is a valid kernel if A is symmetric
 $\Rightarrow A = Q \Lambda Q^T$ where $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$ and $\lambda_1, \lambda_2, \dots, \lambda_n \geq 0$

define $F = \text{diag}(\sqrt{\lambda_1}, \sqrt{\lambda_2}, \dots, \sqrt{\lambda_n})$ then $\Lambda = F^T F$

$$\Rightarrow A = Q^T F^T F Q$$

$$\begin{aligned}\Rightarrow x^T A y &= x^T Q^T F^T F Q y \\ &= (F^T x)^T (F^T y) \\ &= \phi(x)^T \phi(y)\end{aligned}$$

$$\text{where } \phi(x) = F^T Q x$$

Aside:
 $(AB)^T = B^T A^T$
 $(ABC)^T = C^T B^T A^T$

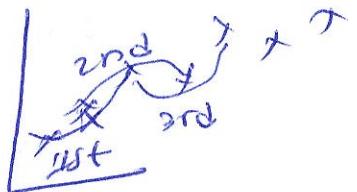
$$\begin{aligned}((AB)C)^T &= C^T (AB)^T \\ &= C^T B^T A^T\end{aligned}$$

⑥ cubic spline smoothing

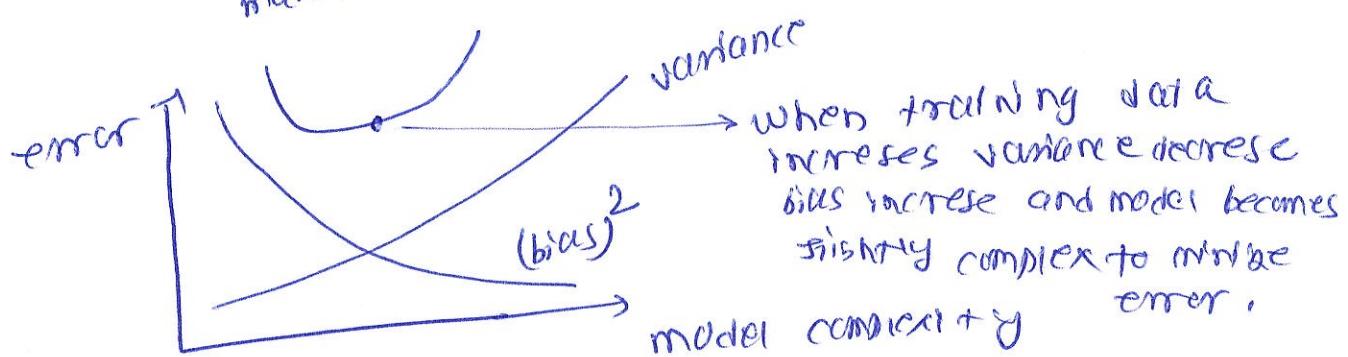
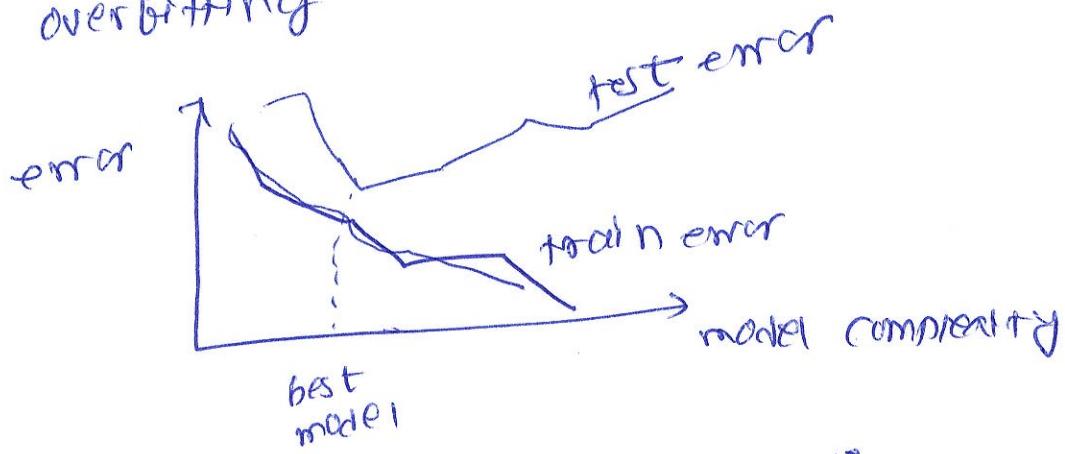
Take 3 points x_i, x_{i+1}, x_{i+2}

$$s_i(x) = a_i (x - x_i)^3 + b_i (x - x_i)^2 + c_i (x - x_i) + d_i$$

$$+ x \in [x_i, x_{i+1}]$$



⑦ overfitting



① Likelihood function for logistic regression

$$P(t_n|w) = b \cdot h_n \cdot (1-h_n)^{1-t_n}$$

$$P(t|w) = \prod_{n=1}^N h_n^{t_n} (1-h_n)^{1-t_n}$$

$$\text{cost } E = \frac{1}{N} \cdot (-\ln P) = -\frac{1}{N} \sum_{n=1}^N [t_n \ln h_n + (1-t_n) \ln (1-h_n)]$$

$$\text{grad} = \nabla J = \frac{1}{N} \cdot \sum_{n=1}^N$$

② perceptron criterion

$$\text{minimize } E_p(w) = -\sum_{n \in \text{Misclassified}} t_n w^T x_n$$

initialize $\bar{w} = 0$, $r = 1$
for $i = 1 \dots N$

$$h_n = \text{sgn}(w^T x_n)$$

if $h_n \neq t_n$

$$w = w + t_n x_n$$

$$r = r + 1$$

$$\bar{w} = \bar{w} + w$$

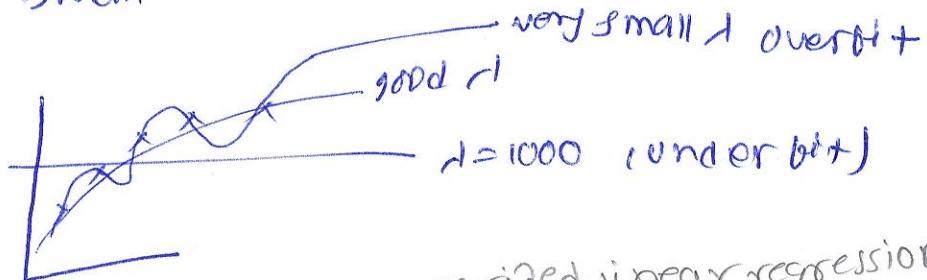
$$r = r + 1$$

return \bar{w}^T

$$\text{Test: } \text{sgn}(w^T x)$$

① L2 \rightarrow more penalty larger weights
does not drive weights to zero
weight vector is NOT sparse.

② A T \rightarrow high bias, low variance, underfit, simpler model
drives weights closer to zero.
small change in training data \Rightarrow big change in estimate



cost function for L2-regularized linear regression

$$③ J = \frac{1}{2N} \sum_{n=1}^N (y_n - t_n)^2 + \frac{\lambda}{2} \|w\|^2$$

$$J(w) = \frac{1}{2N} \sum_{n=1}^N \left(w_0 + \sum_{j=1}^M w_j x_{jn} - t_n \right)^2 + \frac{\lambda}{2} \sum_{j=1}^M w_j^2$$

regularizing w_0 = shifting origin of target

\Rightarrow some change in all target values $\not\Rightarrow$ similar change in estimates

N-examples
M-features

$$X = \begin{bmatrix} x_1 & x_1^2 & x_1^3 & x_1^M \\ x_2 & x_2^2 & x_2^3 & x_2^M \\ x_N & x_N^2 & x_N^3 & x_N^M \end{bmatrix}$$

(we omit w_0 in regularization)

④ Batch GD

for numbers:

$$w^{t+1} = w^t - \eta \nabla J$$

learning rate

stochastic GD

for numbers:

for sample in data:

$$w^{t+1} = w^t - \eta \nabla J$$

$\eta \approx 0.9$

momentum $\mathcal{G} \leftarrow D + m \mathcal{G}$

$$\nabla v^{t+1} = \eta \nabla J(w^t)$$

$$w^{t+1} = w^t - \mathcal{G}$$

$$\nabla v^{t+1}$$

Nesterov accelerated gradient

$$w^{t+1} = w^t - \eta \nabla J(\tilde{w}^t - \eta w^t) - \eta w^t$$

(5) Gradient descent

$\eta = \text{learning rate}$
 $\gamma = \text{momentum}$
 $\eta = \text{hyperparameters}$
 $w^t = \text{weight vector at time stamp } t$

$$v^{t+1} = \begin{cases} \eta \nabla J(w^t) & \text{vanilla} \\ \eta \nabla J(w^t) + \gamma v^t & \text{with momentum} \\ \eta \nabla J(w^t - \gamma v^t) + \gamma v^t & \text{Nesterov} \end{cases}$$

$$w^{t+1} = w^t - v^{t+1}$$

$\gamma \approx 0.9$ Accelerated gradient NAG

(6) Types of GD based on data used:

Batch GD \rightarrow use all data

Stochastic GD \rightarrow use only one example and update w after each iteration

Minibatch GD \rightarrow use constant number N examples and update w after each iteration

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

$$\nabla J = \frac{1}{N} \sum_{n=1}^N (h_n - t_n) \cdot \nabla h_n = \frac{1}{N} \cdot \text{np.sum}((h - t) \cdot X)$$

$$w = w - \eta \text{grad}$$

(7) GD vs Normal eqn

↓

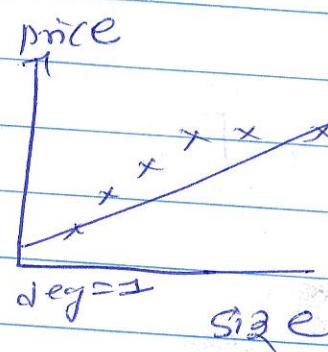
$$w = w - \eta \text{grad}$$

$$w = \underbrace{(X^T X)^{-1}}_{\text{pseudo inverse}} X^T t$$

$$w = \underbrace{(X^T X + \alpha I)^{-1}}_{\text{for } I \text{ identity matrix of shape } X^T X} X^T t$$

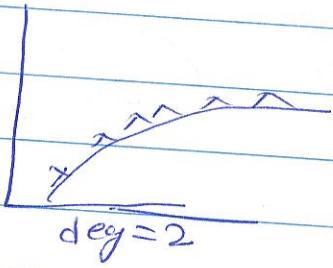
for regularized I is identity matrix of shape $X^T X$
 $I[0,0]=0$ for not to regularize bias term

Bias-variance Trade off



$$w_0 + w_1 x$$

high bias
(underfit)
(simple model)



$$w_0 + w_1 x + w_2 x^2$$

just right

$$w_0 + w_1 x + w_2 x^2 + \dots + w_n x^n$$

high var.

overfit

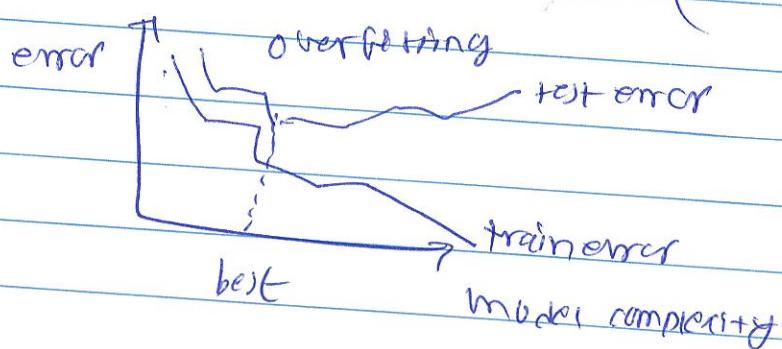
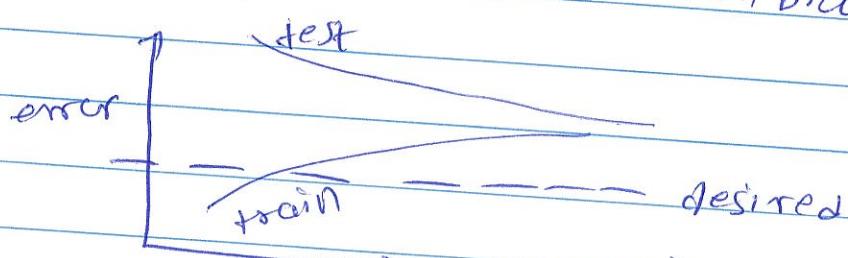
more training eg
smaller feature set
larger feature set

\rightarrow fix high σ^2 (overfitting)

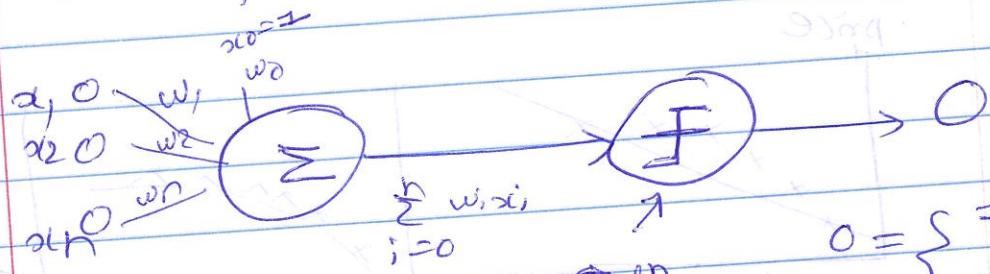
\rightarrow fix high σ^2 (II)

\rightarrow fix high bias
(underfitting)

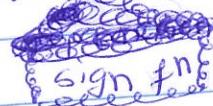
learning curve for high bias (underfit)



perceptron



sign fn



$$O = \begin{cases} 1 & \text{if } \sum_{i=0}^n w_i x_i > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$O(b_1, x_1, x_2, \dots, x_n) = \begin{cases} 1 & \text{if } w_0 + w_1 x_1 + \dots + w_n x_n > 0 \\ 0 & \text{otherwise} \end{cases}$$

(using max)

$$= \begin{cases} 1 & \text{if } \vec{w} \cdot \vec{x} > 0 \\ 0 & \text{otherwise} \end{cases}$$

standardized not sum weighted

(giving constant)

standardized not sum weighted

standardized not sum weighted

① cost for multivariate linear regression

$$\begin{aligned} J &= \frac{1}{2N} \sum_{n=1}^N (h_n - t_n)^2 \\ J &= \frac{1}{2N} \sum_{n=1}^N \left(\sum_{j=0}^M w_j x_n^j - t_n \right)^2 \end{aligned}$$

$$x_1 \quad \begin{bmatrix} 1 & 10 & 100 \\ 1 & 20 & 200 \\ 1 & 30 & 300 \\ 1 & 50 & 500 \end{bmatrix}_{50 \times 4}$$

$$t \quad \begin{bmatrix} 112 \\ 212 \\ 312 \\ 50 \end{bmatrix}_{50 \times 1}$$

2 features ~~age, floorsize~~

$$h = \omega^T X \quad h = x_1 @ \omega^T$$

$$h_n = \sum_{j=0}^M w_j x_n^j$$

$h_1 = w_0 + w_1 x_1^{(1)} + w_2 x_1^{(2)} + \dots + w_M x_1^{(M)}$

$$\omega = [w_0 \ w_1 \ w_2 \ w_3]_{1 \times 4}$$

bias

$$h = (50, 1) \cdot (4, 1) = (20, 1)$$

$$e = h - t = \begin{bmatrix} h_1 - t_1 \\ h_2 - t_2 \\ \vdots \\ h_N - t_N \end{bmatrix}_{50 \times 1}$$

$$SSE = (h - t) \times 2$$

$$MSE = \frac{(h - t) \times 2}{N}$$

$$MSSE = np \cdot \text{mean}(SSE)$$

$$RMSE = np \cdot \sqrt{\text{MSSE}}$$

$$J = \frac{1}{2} \times MSE \quad (\text{single float})$$

$$J = np \cdot \text{sum}((h - t) \times 2) / 2/N$$

$$J = 0.5 \times np \cdot \text{mean}(\sum (h - t) \times 2)$$

$$J = \frac{0.5}{\text{len}(t)} (h - t)^T (h - t)$$

$$(h - t)^T (h - t) = [h_1 - t_1 \ h_2 - t_2 \ \dots]_{1 \times 50} \begin{bmatrix} h_1 - t_1 \\ h_2 - t_2 \\ \vdots \\ h_N - t_N \end{bmatrix}_{50 \times 1}$$

② gradient of J

$$\nabla_{\omega} J = \nabla_{\omega} \frac{1}{2N} \sum_n (h - t)^2 = \left[\frac{\partial J}{\partial \omega_0} \ \frac{\partial J}{\partial \omega_1} \ \frac{\partial J}{\partial \omega_2} \ \dots \ \frac{\partial J}{\partial \omega_M} \right]_{4 \times 1}$$

$$= (h - t) \cdot T @ X^T / N$$

$$= (4, 50) \cdot (50, 1)$$

$$\text{grad} = \nabla_{\omega} J = (1, 4)$$

$$\text{grad-ols} = (h - t) \cdot T @ X^T$$

$$\begin{aligned} \frac{\partial J}{\partial \omega_j} &= \frac{1}{2N} \sum_{n=1}^N (h_n - t_n)^2 = \frac{1}{2N} \sum_{n=1}^N (x_n \omega - t)^T x_n \\ &= \frac{1}{2N} \sum_{n=1}^N 2(x_n \omega - t) \cdot x_n \\ &= \frac{1}{N} \sum_{n=1}^N (x_n \omega - t) \cdot x_n \quad (\text{summation vanishes}) \quad (\text{we use matrix}) \end{aligned}$$

$$\nabla_{\omega} J = \frac{1}{N} (h - t) \cdot X$$

④ $X = \text{design matrix } []_{N \times M}$ $N = \text{samples}$
 $x_1 = \text{biased design matrix } [1]_{N \times M+1}$ $m = \text{features}$

$t = []_{N \times 1}$ column vector

$w = [w_0 \ w_1 \ w_2 \ \dots \ w_m]_{1, M+1}$ row vector (2d array)

$w = np.array(w).rshape(1, \frac{x_1.shape[1]}{1}) = (1, M+1)$

$h = x_1 @ w.T = (N, 1) (M+1, 1) = (N, 1)$ same as t

$$e = h - t = N, 1 = SD, 1$$

linear regression

$$SSE = \sum_{n=1}^N (h - t)^2$$

$$MSE = \frac{1}{N} \sum_{n=1}^N (h - t)^2$$

$$J = \frac{1}{2N} \sum_{n=1}^N (h - t)^2 = np.sum((h - t) ** 2) / 2 / float(N)$$

$$RMSE = E_{RMS} = \sqrt{MSE}$$

Normal eqn: $w = \underbrace{(X^T X)^{-1}}_{\text{matrix penrose pseudo inverse of } X} X^T e + t$
 $\text{np.linalg.pinv}(X)$

Multivariate Linear

bias	x	bed	t
florsize	2	2K	
100	2	3K	
200	3	5K	
500	5		$t \downarrow$

m features and one bias

$$X = \begin{bmatrix} x_0^{(1)} & x_1^{(1)} & \dots & x_m^{(1)} \\ x_0^{(2)} & x_1^{(2)} & \dots & x_m^{(2)} \\ \vdots & \vdots & & \vdots \\ x_0^{(N)} & x_1^{(N)} & \dots & x_m^{(N)} \end{bmatrix}_{N \times m+1}$$

$$t = \begin{bmatrix} t_1 \\ t_2 \\ t_3 \\ \vdots \\ t_n \end{bmatrix}_{n \times 1}$$

→ there are m features and one bias total m+1
 → there are N samples for each feature

$$J = \frac{1}{2N} \|xw - t\|^2 = \frac{1}{2N} (xw - t)^T (xw - t)$$

$$\nabla_w J = \frac{1}{N} \underbrace{\frac{d}{dx} (xw - t)}_{\text{is a matrix}} \cdot \cancel{x^T (xw - t)} \quad \frac{d}{dx} (x^T x) = 2x$$

$$0 = x^T (xw - t)$$

$$x^T x w = x^T t$$

weights

$$w = (x^T x)^{-1} (x^T t)$$

$$\boxed{\frac{d}{dx} (A \cancel{x} + b)^T (A \cancel{x} + b)} \\ = 2A^T (A \cancel{x} + b)$$

derivatives of
products of
matrices

L₁ vs L₂ norm

more used

L₂ → more penalty on large weights, but doesn't drive small weights to zero.

L₁ → less penalty for large wt, but tends many weights towards (or very very close) to zero.
leading to weight vector to be sparse.

Bias-variance Tradeoff

$$\textcircled{a} \text{ Ridge: } \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij})^2 + \lambda \sum_{j=1}^p \beta_j^2 = RSS + \lambda \sum_{j=1}^p \beta_j^2$$

shrinkage parameter
which has to be
tuned via validation
set

→ regularization
strength

ridge meaning:
 ↗ small change in training data
 ↘ big change in parameter estimates
 effect will increase with no. of parameters

$$\textcircled{b} \text{ Lasso } \beta_j^2 \leq |\beta_j|$$

Ridge shrinks wj but don't make them 0
 but lasso makes them 0, makes less non-zero

Logistic Regression \rightarrow classification (binary)

linear regression $h = \omega^T x$

logistic regression $h = \sigma(\omega^T x) = \frac{1}{1 + e^{-\omega^T x}} = h(x) = y(b)$

likelihood function $p(t|w) = \prod_n h_n^{t_n} (1 - h_n)^{1-t_n}$

negative log likelihood, E or $J = -\ln p(t|w)$

$$= -\sum_{n=1}^N [t_n \ln h_n + (1-t_n) \ln (1-h_n)]$$

i. cost or error or loss function

$$E_D = -\frac{1}{N} \sum [t_n \ln h_n + (1-t_n) \ln (1-h_n)]$$

where $t_n \in \{0, 1\}$ NOT $\{-1, 1\}$

add regularization

$$E_w = \frac{\lambda}{2} \omega^T \omega$$

then L2 regularized logistic regression cost function

$$\text{cost} \# E = -\frac{1}{N} \sum [t_n \ln h_n + (1-t_n) \ln (1-h_n)] + \frac{\lambda}{2} \omega^T \omega$$

$$h(x) = \sigma(\omega^T x) = \frac{1}{1 + e^{-\omega^T x}}$$

$$h_n = h(x(n)) = \sigma(\omega^T x(n)) = \frac{1}{1 + e^{-\omega^T x(n)}} \text{ is sigmoid fn}$$

softmax regression

training set: $(x_1, t_1), (x_2, t_2), \dots, (x_n, t_n)$

$$x = [1, x_1, x_2, \dots, x_n]$$

$$t_1, t_2, \dots, t_n \in \{1, 2, \dots, K\}$$

$$w_{1K} = [w_{10}, w_{11}, w_{12}, \dots]^T$$

one weight vector per class,

$$p(t_k|x) = \frac{e^{w_k^T x}}{\sum_j e^{w_j^T x}}$$

prob that
n-th example x_n
has class t_n

$$p(t_n|x_n) = \frac{e^{w_{tn}^T x_n}}{\sum_j e^{w_j^T x_n}}$$

likelihood is the joint probability of all classes

$$L(w) = \prod_{n=1}^N p(t_n|x_n)$$

cost is the -ve log likelihood,

$$E_D(w) = -\frac{1}{N} \ln L(w)$$

$$= -\frac{1}{N} \ln \prod_{n=1}^N p(t_n|x_n)$$

$$= -\frac{1}{N} \sum_n \ln p(t_n|x_n)$$

$$= -\frac{1}{N} \sum_n \ln \frac{e^{w_{tn}^T x_n}}{\sum_k p(t_k)}$$

$$E_D(w) = -\frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K \delta_{k,t_n} \ln \left(\frac{e^{w_k^T x_n}}{\sum_k p(t_k)} \right)$$

MLE

sample: x_1, x_2, \dots, x_N

probability of x is $b(x, w)$

joint prob of x_1, x_2, \dots, x_N is $\prod b(x_n, w)$

Likelihood $L(w)$

$x_N = x_n$

$L(w) = p(x=x_n | w)$

$= b(x_n, w)$

$\prod b(x_n, w)$

$L(w) = \prod_{n=1}^N b(x_n, w)$

MAP solution for LR

$$p(t|w) = \prod_{n=1}^N h_n^{t_n} (1-h_n)^{1-t_n}$$

$$p(w) = (\frac{2}{\pi})^{\frac{m+1}{2}} e^{-\frac{d}{2} w^T w} \propto e^{-\frac{d}{2} w^T w}$$

$$p(w|t) = \frac{p(t|w) p(w)}{p(t)} \propto p(t|w) p(w)$$

$$\text{now, } \hat{w}_{\text{MAP}} = \underset{w}{\operatorname{argmax}} \ p(w|t)$$

$$= \underset{w}{\operatorname{argmax}} \ p(t|w) p(w)$$

$$= \underset{w}{\operatorname{argmin}} -\ln p(t|w) -\ln p(w)$$

$$= \underset{w}{\operatorname{argmin}} -\ln \prod_n h_n^{t_n} (1-h_n)^{1-t_n} -\ln e^{-\frac{d}{2} w^T w}$$

$$= \underset{w}{\operatorname{argmin}} -\sum_n [t_n \ln h_n + (1-t_n) \ln (1-h_n)] + \frac{d}{2} w^T w$$

$$= \underset{w}{\operatorname{argmin}} \underbrace{-\sum_n [t_n \ln h_n + (1-t_n) \ln (1-h_n)]}_{E_D(w)} + \underbrace{\frac{d}{2} w^T w}_{E_W(w)}$$

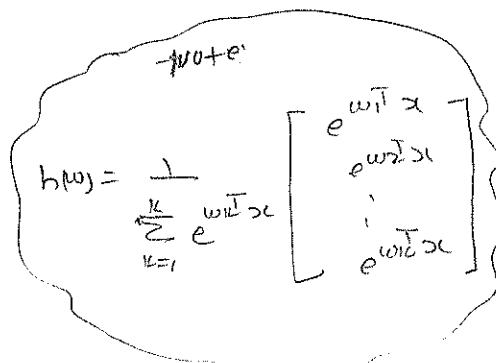
$$\boxed{\hat{w}_{\text{MAP}}} = \underset{w}{\operatorname{argmin}} E_D(w) + E_W(w)$$

ans

softmax Regression

$$P(c_k|x_n) = \frac{e^{w_k^T x_n}}{\sum_j e^{w_j^T x_n}}$$

$$P(t_n|x_n) = \frac{e^{w_{t_n}^T x_n}}{\sum_{j=1}^K e^{w_j^T x_n}}$$



$$E_D(w) = -\frac{1}{N} \sum_n \prod_i P(t_n|x_n)$$

$$= -\frac{1}{N} \sum_n \prod_i \frac{e^{w_i^T x_n}}{\sum_j e^{w_j^T x_n}}$$

$$= -\frac{1}{N} \sum_n \prod_i \frac{e^{w_i^T x_n}}{\sum_j e^{w_j^T x_n}}$$

regularizer

$$E_D(w) = -\frac{1}{N} \sum_n \sum_k \delta_{ik}(t_n) \ln \left(\frac{e^{w_k^T x_n}}{\sum_l e^{w_l^T x_n}} \right)$$

$$+ \frac{\alpha}{2} \sum_{k=1}^K w_k^T w_k$$

$$\frac{\partial E_D}{\partial w_k} = \alpha w_k$$

(no summation)

$$\text{Let, } E_n = \sum_k \delta_{ik}(t_n) w_k^T x_n - \sum_k \delta_{ik}(t_n) \ln \left(\sum_k e^{w_k^T x_n} \right)$$

$$\frac{\partial E_n}{\partial w_j} = \delta_{ij}(t_n) x_n - \delta_{ij}(t_n) \cdot \frac{1}{\sum_j e^{w_j^T x_n}} \cdot e^{w_j^T x_n} \cdot \delta_{ij}(t_n) x_n$$

$$\frac{\partial E_n}{\partial w_k} = \delta_{ik}(t_n) x_n - \frac{e^{w_k^T x_n}}{\sum_l e^{w_l^T x_n}} \cdot x_n$$

prevent overfitting

$$p(c_k|x_n) = \frac{e^{w_k^T x_n}}{\sum_l e^{w_l^T x_n}}$$

$$c = \max_{1 \leq k \leq K} w_k^T x_n$$

$$p(c_k|x_n) = \frac{e^{(w_k^T x_n - c)}}{\sum_l e^{(w_l^T x_n - c)}}$$

$$\frac{\partial E_n}{\partial w_k} = [\delta_{ik}(t_n) - p(c_k|x_n)] x_n$$

$$\Rightarrow \frac{\partial E_D(w)}{\partial w_k} = -\frac{1}{N} \sum_n [\delta_{ik}(t_n) - p(c_k|x_n)] x_n$$

there are K such equations for each classes.

① MLE VS MAP

MLE = maximize $P(\text{data} | \text{params})$ by searching over parameters

MAP = maximize $P(\text{param} | \text{data})$ by searching over params and accounting for prior over params.

MLE \rightarrow finds ω by maximizing likelihood $P(\text{data} | \omega)$
MAP \rightarrow minimizes the posterior prob $p(\omega | \text{data})$

② classification maps inputs to discrete outputs
Regression maps inputs to continuous outputs

③ PCA & feature selection
similarity : reduce the dimension of data
difference : feature selection finds a subset of features
PCA produces a smaller new set

perceptron

perceptron criterion: $w^T x_n > 0$ for $t_n = +1$
 $w^T x_n < 0$ for $t_n = -1$

want: $t_n w^T x_n > 0$ for all patterns

minimize $\sum -w^T x_n t_n$ for all misclassified patterns M

$$E_p(w) = \sum_{n \in M} w^T x_n t_n$$

perception mistakes (misclassified)

Binary perceptron

initialize $\vec{w} = 0$

for $n = 1, \dots, N$

$$h_n = \text{sgn}(w^T x_n)$$

if $h_n \neq t_n$ then

$$w = w + t_n x_n$$

} repeat until convergence
or
given number of epochs

① why Kernels are symmetric?

Inner products are symmetric by definitions, so, therefore if the kernel function represent an inner product in some Hilbert space, then the kernel function must be symmetric as well.

$$\begin{aligned}K(x,y) &= \langle \phi(x), \phi(y) \rangle \\&= \langle \phi(x), \phi(y) \rangle \quad (\because \text{property of inner product}) \\&= K(y,x)\end{aligned}$$

② show $\langle \phi(x), \phi(y) \rangle = \alpha^T A^T A y$ is a valid kernel

Let $\phi(x) = Ax$,

then,

$$\begin{aligned}\langle \phi(x), \phi(y) \rangle &= \phi(x)^T \phi(y) \\&= (Ax)^T (Ay) \\&= \alpha^T A^T A y \\&= K(x,y)\end{aligned}$$

$\therefore \langle \phi(x), \phi(y) \rangle$ is an inner product in some Hilbert Space.

③ $A^T A$ is PSD matrix (indirect)

Proof Let, $A \in \mathbb{R}^{m \times n}$

λ = eigen value of $A^T A$

q = eigen vector of λ

$$(A^T A)q = \lambda q \quad (\text{premultiplied by } q^T)$$

$$\Rightarrow q^T A^T A q = q^T \lambda q$$

$$\therefore A^T A = A^T A$$

$$\begin{aligned} \lambda &= \frac{q^T A^T A q}{q^T q} = \cancel{\frac{q^T A^T A q}{q^T q}} \\ &= (q^T A^T A q) / q^T q \end{aligned}$$

$$= \frac{z^T z}{q^T q} \quad \text{where } z = A^T q$$

$$\text{here } z^T z \geq 0$$

$$q^T q \geq 0$$

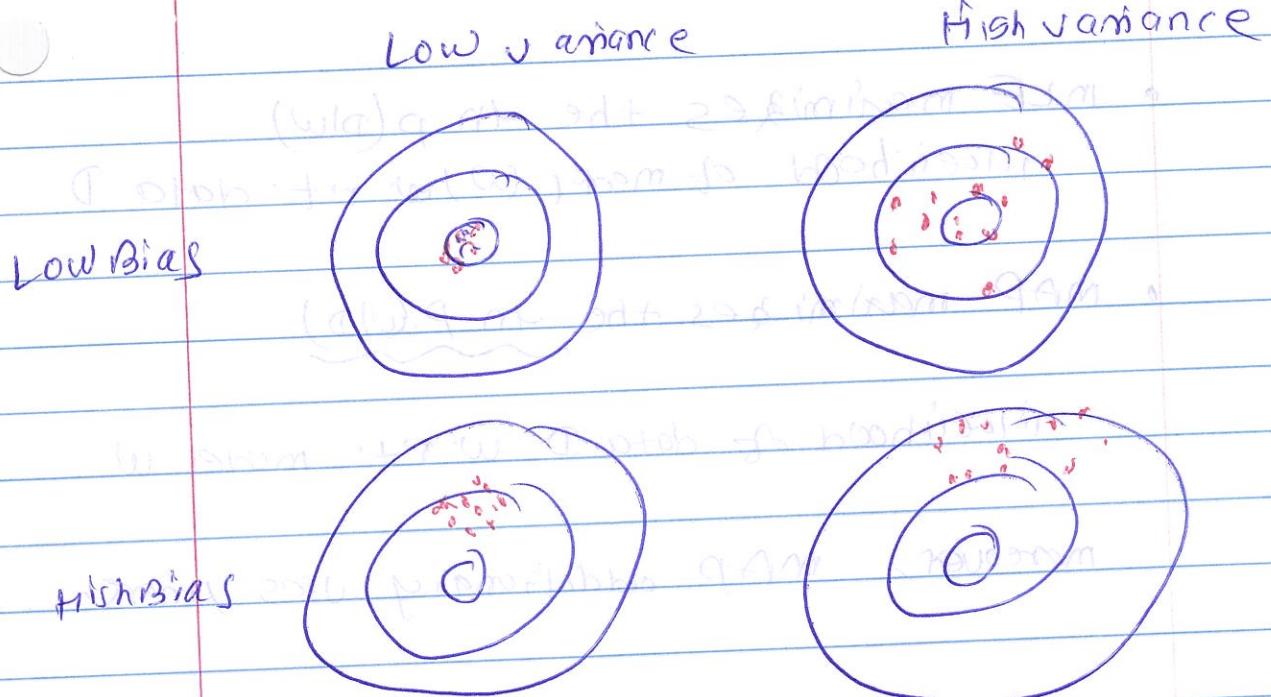
$\therefore \lambda \geq 0 \text{ so } A^T A \text{ is PSD.}$

MLE v/s MAP

- MLE maximizes the $-\ln p(D|w)$
likelihood of model w wrt. data D
 - MAP maximizes the $\ln P(w|D)$
likelihood of data D wrt. model w
- moreover, MAP additionally uses priors.

QAM & V.EJM

Q.



Error

optimum
model/
complexity

total error

variance

$Bias^2$

model complexity

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

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

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

$$K(x_1, x_2) = (x_1^T x_2)^2 = (x_1^1 x_2^1 + x_1^2 x_2^2)^2$$

$$= (x_1^1)^2 + 2 x_1^1 x_1^2 x_2^1 x_2^2 + (x_1^2)^2$$

$$x = \begin{pmatrix} x_1^1 \\ x_1^2 \end{pmatrix}$$

$$x^T = \begin{pmatrix} x_1^1 & x_1^2 \end{pmatrix}$$

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

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

thus, gives, $\frac{q(x)}{\phi(x) \phi(x)} = \frac{x_1^2 + x_2^2 + 2 x_1 x_2}{x_1^2 + 2 x_1 x_2 + x_2^2}$

$$K(x_1, x_2) = \phi(x_1)^T \phi(x_2)$$

* A necessary and sufficient condition for a kernel function to be "valid" is that the Gram matrix be positive and semi-definite for all choices of $\phi(x)$.

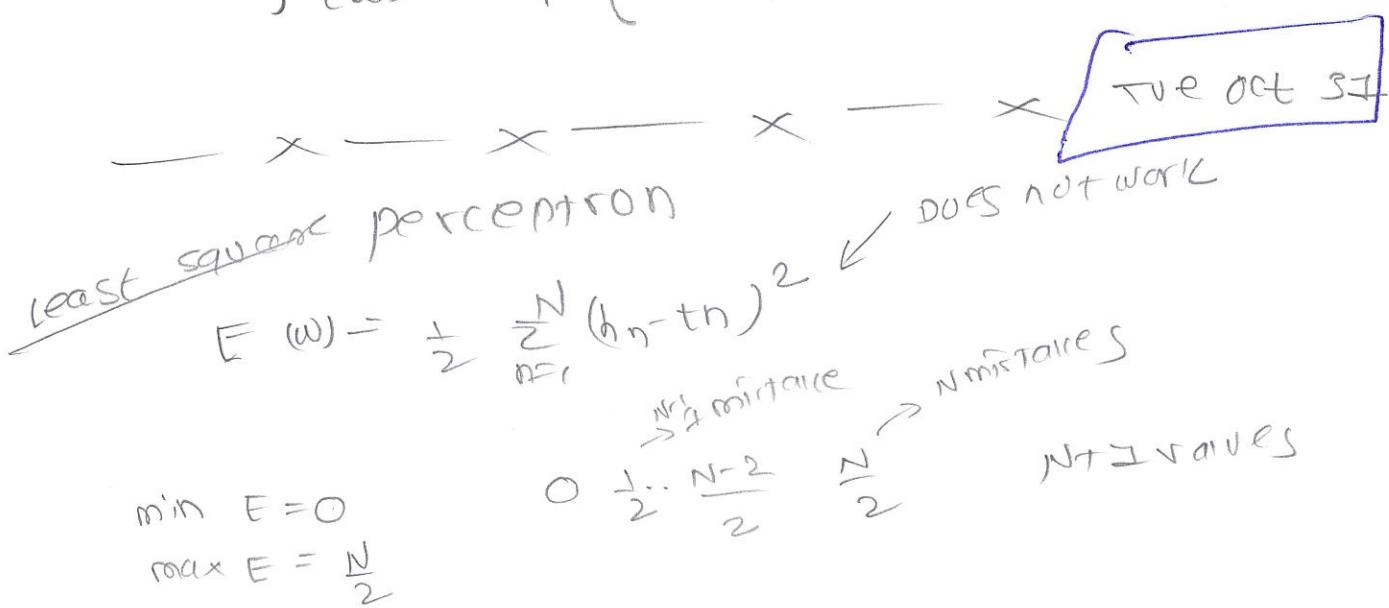
A Gram matrix of $\phi(x)$ is $x^T x$.

The linear vector x is projected onto a quadratic surface.

If all the points in this surface are non-zero, then our kernel is valid.

* Fisher's discriminant cost

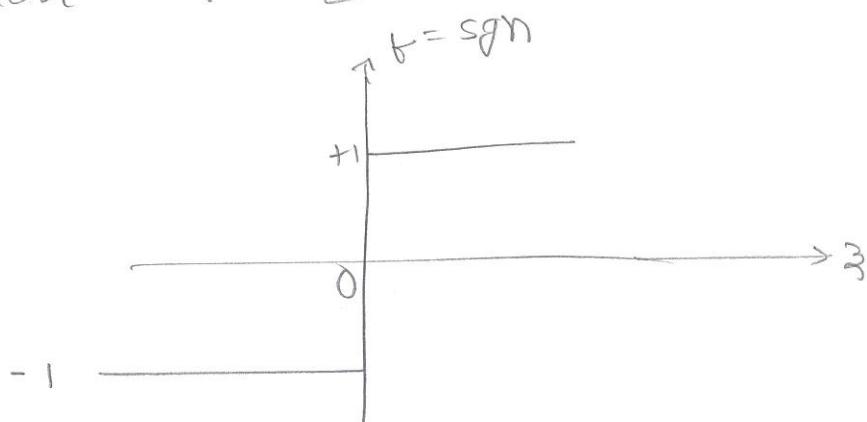
$$J(w) = \text{Tr} \{ w s_B w^T \} \cdot (w s_B w^T)$$



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

cost = no. of misclassified patterns
= 0 for no misclassification
 $\frac{N}{2}$ for N mistakes

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



Distance metrics

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

② Hamming $d(x, y) = \# \text{ of different values in binary length strings}$

③ Mahalanobis distance $d(x, y) = \sqrt{(x-y)^T S^{-1} (x-y)}$
 S is sample cov. matrix

$S = I \rightarrow \text{Euclidean}$

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

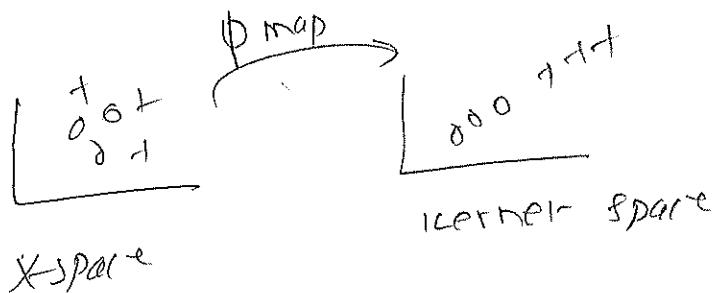
④ cosine similarity

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

⑤ Levenshtein distance (edit distance)

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

$$x = \text{'attens'} \quad y = \text{'hnts'} \quad d(x, y) = 4$$



~~①~~ SVM with slack (soft margin SVM)

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

with constraint

$$t_n (w^\top \phi(x_n) + b) \geq 1 - \xi_n$$

$$\xi_n \geq 0 \quad \text{and} \quad \sum_{i=1}^N \xi_i \leq \varepsilon \quad \forall n \in \{1, 2, \dots, N\}$$

i.e. $1 - \xi_n - t_n w^\top \phi(x_n) - b t_n \leq 0 \quad \text{--- ① (constraint)}$

$-\xi_n \leq 0 \quad \text{--- ② (constraint)}$
 $\forall n = 1, 2, \dots, N$

primal Lagrangian

$$L_p(w, b, \xi, \alpha, r) = \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N \xi_n + \sum_{n=1}^N \alpha_n (1 - \xi_n - t_n w^\top \phi(x_n) - b t_n) - \sum_{n=1}^N r_n \xi_n$$

dual Lagrangian

$$L_d(\alpha, r) = \inf_{w, b, \xi} L_p(w, b, \xi, \alpha, r)$$

Note / $\delta = \frac{\partial L_p(w, b, \xi, \alpha, r)}{\partial w} = w - \sum_n \alpha_n t_n \phi(x_n) \Rightarrow (w = \sum_n \alpha_n t_n \phi(x_n))$

$$\delta = \frac{\partial L_p}{\partial b} = - \sum_n \alpha_n t_n \Rightarrow (\sum_n \alpha_n t_n = 0)$$

NO summation! $\delta = \frac{\partial L_p}{\partial \xi_n} = C - \alpha_n - r_n \Rightarrow (C = \alpha_n + r_n) \quad \forall n \in \{1, 2, \dots, N\}$

then,

$$L_D(K) = \frac{1}{2} \|w\|^2 + C \sum_n \xi_n + \sum_n \alpha_n (1 - \xi_n - t_n w^\top \phi_n - t_n b) - \sum_n r_n \xi_n$$

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

then the optimization problem in dual space is,

$$\underset{\alpha}{\text{maximize}} \quad L_D(K), = \sum_n \alpha_n - \frac{1}{2} \sum_{m,n} \alpha_m \alpha_n t_m t_n K^{(m,n)}$$

with constraints $0 \leq \alpha_n \leq C \quad \forall n = 1, 2, \dots, N$

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

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

$$\text{since, } C \sum_n \xi_n = \sum_n \alpha_n \xi_n + \sum_n r_n \xi_n \\ = \sum_n (t_n + r_n) \xi_n$$

- ①
- | | | |
|---|---|---------------------------------|
| 0 | + | • SVM cannot kernel & can occur |
| + | 0 | • zero training error |
- Logistic Regr, 3-NN cannot

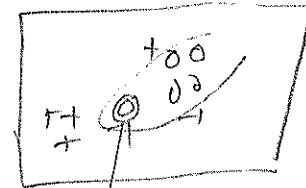
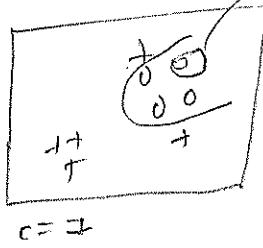
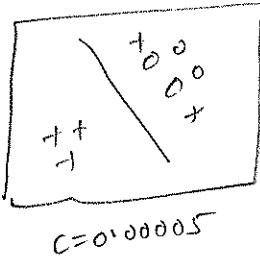
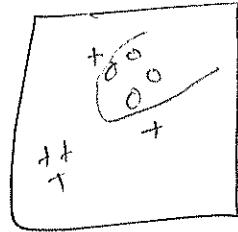
- ② If examples are iid,
increase training examples
- may increase train error
but decrease test error

- ③ SVM effect of C

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

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

will not change
decision boundary



$C = 10,000$

between $C \approx 1$ 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.

adding this change
decision boundary

- ④ Bias variance Tradeoff

linear reg
 $d=2$ poly
 $d=10$ poly

Bias	Vari
high	$\ w\ $
low	$\ w\ $
low	high

python

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

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

$\text{ones} = np.ones(x.shape[0]).reshape(-1, 1)$

$x_2 = np.append(\text{ones}, x, axis=1).astype(np.int)$

$x_2 = np.c_[np.ones(x.shape[0]).reshape(-1, 1), x]$

$x_2 = np.column_stack((1,))$

$x_2 = np.c_[np.ones(x.shape[0])[:, np.newaxis], x]$

$x_2 = np.c_[np.ones(x.shape[0]).
 \overbrace{\text{reshape}(-1, 1)}$

$= np.c_[np.atleast_2d(np.ones(x.shape[0]).T), x]$

$= np.c_[np.expand_dims(np.ones(x.shape[0], axis=1), 1), x]$

⑤ Given $\phi(x) = [1, x_1, x_2, x_1x_2]$

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

$$\Rightarrow K(x, x') = 1 + x_1x_1' + x_2x_2' + x_1x_2 x_1'x_2'$$

⑥ L1 vs L2 loss

L2 is more robust to outlier than L1

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

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

⑤ L1 gives sparse solution used in feature selection.

⑥ Logistic loss is better than L2 loss in classification.

⑦ SVM small C,

for linearly separable data, small C can affect
training accuracy.

A small C can allow large margins thus, the
resulting classifier will have smaller w^2
and can have non-zero training error.

① solve the SVM problem without slack using Lagrange multiplier method

→ the optimization problem is

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

$$\text{s.t. } t_n (\omega^\top \phi^{(n)} + b) \geq 1 \quad \text{at } n=1, \dots, N$$

$$1 \leq t_n (\omega^\top \phi^{(n)} + b)$$

$$1 \leq t_n \omega^\top \phi^{(n)} + t_n b$$

$$1 - t_n \omega^\top \phi^{(n)} - t_n b \leq 0 \quad \text{(convex constraint)}$$

compare $t_i (a_i) \leq 0$ then $a_i = 1, \dots, m$

primal Lagrangian,

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

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 :

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

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

LOOK LHS

then dual lagrangian is,

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

~~$-\sum_n \alpha_n \text{tr}(\phi_n^2)$~~

$$L_D(d) = \sum_n \alpha_n - \frac{1}{2} \sum_{n,m} \alpha_n \alpha_m \text{tr}(\phi_n \phi_m)$$

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

Then the optimization problem in dual space is,

$$\text{maximize } L_D(d) = \sum_n \alpha_n - \frac{1}{2} \sum_{n,m} \alpha_n \alpha_m \text{tr}(\phi_n \phi_m)$$

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

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

Now, KKT conditions are,

① primal constraints $1 - \alpha_n (\vec{w}^T \phi(x_n) + b) \leq 0 \quad (\text{convex constraint equation})$
 $\alpha_n (\vec{w}^T \phi(x_n) + b) \geq 1$

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

③ complementary slackness $\alpha_n \{1 - \alpha_n (\vec{w}^T \phi(x_n) + b)\} = 0$

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

$$1 - \alpha_n (\vec{w}^T \phi(x_n) + b) = 0$$

these α_n are called support vectors.

Ques here $n = 1, 2, \dots, N$

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.

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

$$\text{or, } 1 - t_m \phi(b_m) \sum_n \alpha_n \phi(b_n) - t_m b = 0$$

$$\text{or, } t_m b = 1 - t_m \phi(b_m) \sum_n \alpha_n \phi(b_n)$$

$$= 1 - t_m \sum_n \alpha_n \phi(x_n) \phi(b_m)$$

$$b_{tm} = 1 - t_m \sum_n \alpha_n \phi(x_n) \phi(b_m)$$

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

$$\boxed{b = t_m - \sum_n \alpha_n \phi(x_n) \phi(b_m)} \quad \because \frac{1}{t_m} = t_m$$

$$\frac{1}{t_m} = \frac{1}{\frac{1}{2}} \text{ and } -1 = \frac{-1}{-1}$$

This is true for all the m examples which have non-zero
Lagrange parameter α .

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

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

where S is the subset of all the examples
where lagrange parameter α is non-zero.

$$S \subseteq \mathbb{D}$$

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

then, Linear discriminant function is,

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

KNN = memory-based
(no model to fit)

① { find k nearest examples t_1, t_2, \dots, t_K from test set.
 $y(t) = \arg \max_{t \in T} \sum_{i=1}^k w_i \delta_T(t_i)$

wi gives distance-weighted KNN
 $w_i = \frac{1}{\|x - t_i\|^2}$

② Mahalanobis distance

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

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

③ cosine similarity

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

④ kernel-based distance weighted NN

→ binary classification $T \in \{-1, 1\}$

$$\rightarrow y(t) = \text{sign} \left(\sum_{i=1}^N k(t, t_i) \cdot T_i \right)$$

- KNN usage
- ① satellite image classification
 - ② digit recognition

Wrapper method

- ① ^{forward} forward selection
- $F = \text{all features}$
- $S = \text{subset of features}$
- start $S = \{\}$ empty
- for each feature b in $F - S$
find best performing because $b \in S$ and add to S .
repeat until performance does not increase
or performance good enough

- ② ^{recursive} backward elimination
- $F = \{1, 2, \dots, K\}$ is set of features
- $S = \{\}$ named set of features
- repeat until $F - S$ is empty
train \vec{w} using linear SVM and $F - S$
find feature b with minimum $|\vec{w}_b|$
append b to S
- return S

⑤ instance-weighted kNN for regression

1. find k nearest points $\mathbf{x}_1 \mathbf{x}_2 \dots \mathbf{x}_k$

$$2. y(\mathbf{x}) = \frac{\sum_{i=1}^k w_i t_i}{\sum_{i=1}^k w_i} \quad \text{where } w_i = \frac{1}{\|\mathbf{x}_i - \mathbf{x}\|^2}$$

$k=N \rightarrow$ stepwise method

$$y(\mathbf{x}) = \frac{\sum_{i=1}^N K(m_i) t_i}{\sum_{i=1}^N K(m_i)} \quad (\text{kernel-based dist-weighted})$$

⑥ Regression with kNN

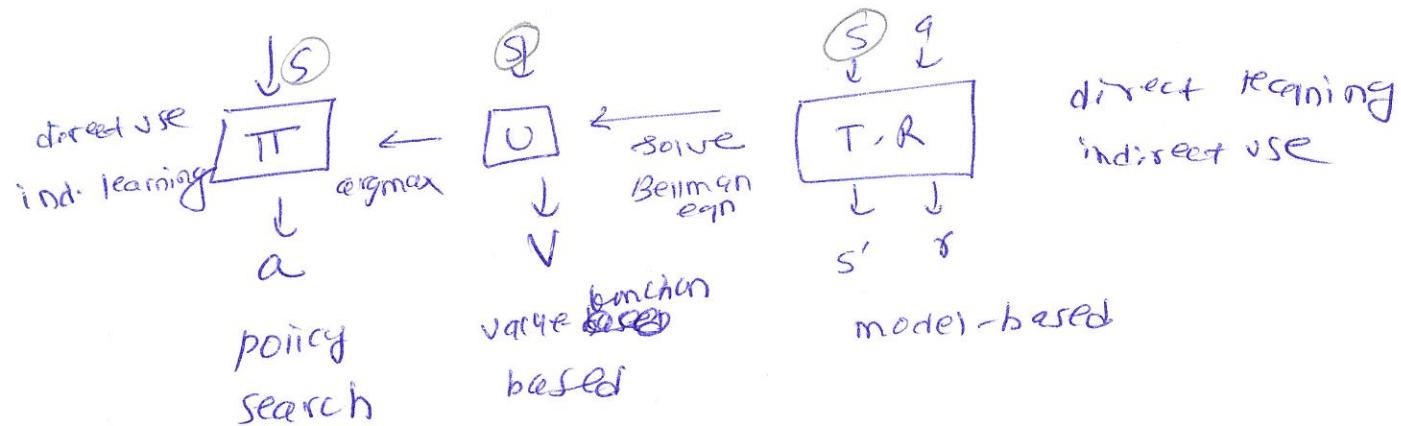
$$y(\mathbf{x}) = +K \sum_{i=1}^k t_i \quad (\text{ti are values, not classes})$$

⑦ distance-weighted kNN (regression)

$$\text{prioral: } y(\mathbf{x}) = \frac{\sum_{i=1}^k w_i t_i}{\sum_{i=1}^k w_i} \quad w_i = \frac{1}{\|\mathbf{x}_i - \mathbf{x}\|^2}$$

$$\text{kernel based: } y(\mathbf{x}) = \frac{\sum_{i=1}^N K(\mathbf{x}, \mathbf{x}_i) t_i}{\sum_{i=1}^N K(\mathbf{x}, \mathbf{x}_i)} \quad (\text{ti are values not classes})$$

3 Approaches to RL



Filter method of Feature Selection

① Mutual Information

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

eq when X, Y indep
max when $x=y$

Let there are K examples with l features.
 $O_{ij} = \begin{cases} \text{observed value} \\ \text{for } x=i, y=j \end{cases}$

$$\begin{matrix} 1 & - & j & - & l \\ \vdots & & O_{ij} & \} & N_{x=i} \\ K & \sim & N_{y=j} & & E_{ij} = \frac{N_{x=i} N_{y=j}}{N} \end{matrix}$$

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

Classification

② probabilistic Generative models

{ naive Bayes

{ hidden markov models

- inf. dec \rightarrow separate

- can use $p(\theta)$ for outlier or novelty

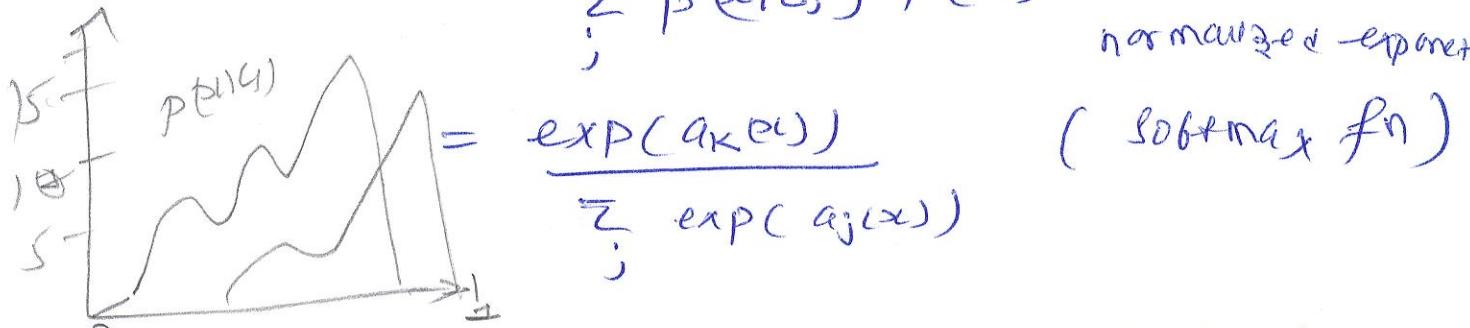
- need to model dependencies between features

Naive Bayes \rightarrow resilient to noise

- text classification with NB \rightarrow OOK HW
- multiple classes posterior prob of class c_k given ~~data~~ x is

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

normalized exponentials



generative where $a_{ik}(x) = \ln p(x|c_k) \cdot b(c_k)$

① SVM for ranking

optimization problem

$$\text{minimize } J(w, b) = \frac{1}{2} \|w\|^2 + C \sum \epsilon_{ik,ij}$$

$$\text{s.t. } w^T \phi(q_k, d_i) \geq w^T \phi(q_k, d_j) + 1 - \epsilon_{ik,ij}$$

$$\epsilon_{ik,ij} \geq 0$$

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

① Add ones column
= $\text{np}\cdot\text{c} \leftarrow \text{np}\cdot\text{ones}(X\cdot\text{shape}[0]) \cdot \text{reshape}(-1, 1), X]$
 $X = \text{np}\cdot\text{c} - \text{c} \leftarrow \text{np}\cdot\text{ones}(X\cdot\text{shape}[0]) \text{ [np.newaxis]} \cdot \text{T}, X]$

② correct = $\text{np}\cdot\text{sum}(y\text{-pred} == y\text{-test})$
accuracy = $\text{correct} / \text{len}(y\text{-pred})$

⑤ Gradient Descent (GD or BGD) GD with momentum

vanilla $v^{t+1} = \eta \nabla J(w^t)$ $v^{t+1} = \gamma v^t + \eta \nabla J(w^t)$
 $w^{t+1} = w^t - v^{t+1}$ $w^{t+1} = w^t - \eta \nabla J(w^t)$

Batch Gradient Descent (Lect 01, p. 24)

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

$$\begin{aligned} w^{t+1} &= w^t - \eta \nabla J(w^t) \\ &= w^t - \eta \sum_{n=1}^N (h_n - t_n) \alpha_n \\ &= w^t - \text{learningRate} (h - t) @ X_1 \end{aligned}$$

confusion matrix

ACTUAL class

		cat	not cat	
predicted class	cat	TP	FP	P
	not cat	FN	TN	N
		true		
		P	N	

$$\begin{aligned} 00 &= TN \\ 01 &= FP \\ 10 &= FN \\ 11 &= TP \end{aligned}$$

		+	0	
1	+	TP	FP	
	0	FN	TN	
		P	N	

		+	0	
0	+	TP	FN	
	0	TN	FP	
		P	N	

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

(hit rate)

$$\frac{\text{predicted true}}{\text{actual true}}$$

F1 score is the harmonic mean

of precision and recall

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

$$= \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

* prove that the number of elements in
x and y is also same!

$\Rightarrow \text{IC}(x,y) = |\text{X} \cap \text{Y}|$ is a kernel.

contd (contd)

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

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

$$+ \sum_{m \notin S} \alpha_m t_m K(x_m, x)$$

$m \notin S \Rightarrow$ but if $x_m \notin S$, then $\alpha_m = 0$

* package libsvm or SUMSVM

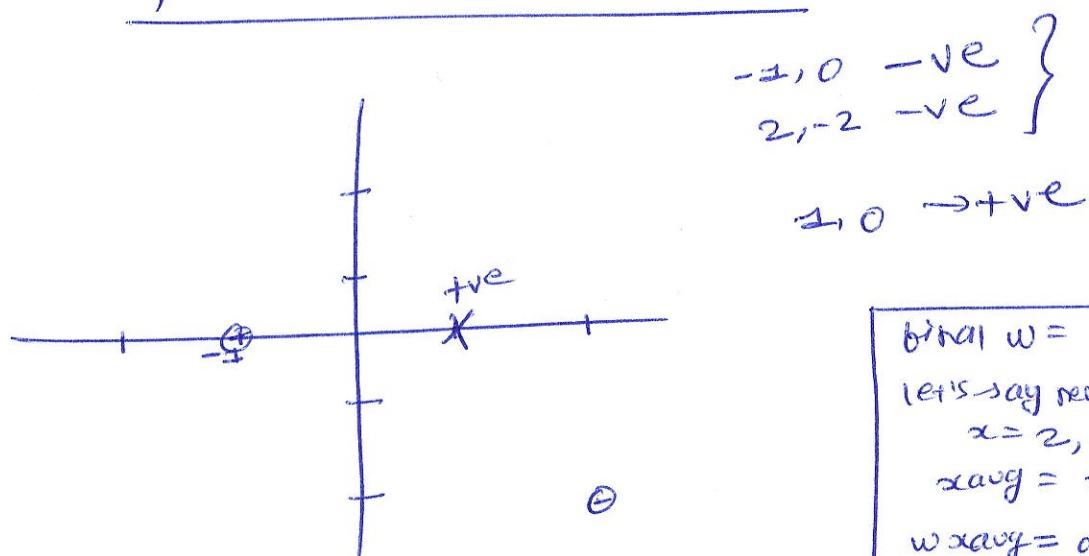
$$\text{input} \rightarrow \{\alpha_n, t_n\} \quad \forall n \in N$$

$$\text{output} \rightarrow \{\alpha_m, t_m, x_m\} \quad m \in S$$

model.txt \rightarrow

$$\boxed{\begin{matrix} b \\ \alpha_m t_m \cup x_m \end{matrix}}$$

perceptron update



final $w = 012$
 let's say new point
 $x = 2, -1.0$
 $x_{avg} = 1.5, -1.0$,
 $w \cdot x_{avg} = 0 + 2 - 2.02 < 0$
 i.e. -ve label

Augmented data

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

X_{avg} t

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

note we can
also use
 $\vec{w} = (0.3, 0.2, 0.1)$
 $b = 0.9$
 original unaugmented x

updated weights correct example

Test point hypothesis correct

eg1 -: (1 -1 0) $w \cdot x = 0 + 0 + 0 \leq 0$ false

$$w = w - x = (-1 \ 1 \ 0)$$

eg2 -: (1 2 -2) $w \cdot x = -1 + 2 + 0 \leq 0$ false

$$w = w - x = (-2 \ -1 \ 2)$$

eg3 +: (1 1 0) $w \cdot x = -2 - 1 + 0 \geq 0$ true

$$w = w + x = (-1 \ 0 \ 2)$$

eg1 -: (1 -1 0) $w \cdot x = -1 + 0 + 0 \leq 0$ true

$$w = w = (-1, 0, 2)$$

eg2 -: (1 2 -2) $w \cdot x = -1 + 0 + 4 \geq 0$ true

$$w = w + x = (0 \ 1 \ 2)$$

eg3 +: (1 1 0) $w \cdot x = -1 + 0 + 0 \geq 0$ true

$$w = w = (0 \ 1 \ 2)$$

Final weight

eg1 -: (1 -1 0)

eg2 -: (1 2 -2)

eg3 +: (1 1 0)

$w \cdot x = 0 + (-1) + 0 \leq 0$ true

$w \cdot x = 0 + 2 + (-4) \leq 0$ true

$w \cdot x = 0 + 1 + 0 \geq 0$ true

① constrained optimization

(Lagrange multipliers)

$$\text{maximize } S = (x_1 - 2)^2 + 2(x_2 - 1)^2$$

$$\text{s.t. } x_1 + 4x_2 = 3$$

soln: If we ignore constraint we get $x_1 = 2, x_2 = 1$
 then $x_1 + 4x_2 = 2 + 4 \cdot 1 = 6$ is too large
 for the constraint.

consider

$$L = L(x_1, x_2, \lambda) = S - \lambda(x_1 - 2)^2 - \lambda(x_2 - 1)^2$$

look at:

$\lambda = 0$ $\lambda = 1$ $\lambda = \frac{2}{3}$	$2, 1$ $3/2, 0$ $\left(\frac{5}{3}, \frac{1}{3}\right)$	$\frac{5}{3} + 4 \cdot \frac{1}{3} = \frac{5}{3} + \frac{4}{3} = \frac{9}{3} = 3$	\checkmark
---	---	---	--------------

formal soln

$$\frac{\partial L}{\partial x_1} = -2(x_1 - 2) - \lambda = 0$$

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

$$2x_1 = 4 - \lambda$$

$$= 4 - \frac{2}{3}$$

$$= 4 - \frac{2}{3} = \frac{10}{3}$$

$$x_1 = \frac{5}{3}$$

~~so~~

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

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

$$\Rightarrow x_2 = 1 - \frac{\lambda}{4}$$

$$x_2 = 1 - \lambda$$

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

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

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

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

$$-6 + 9\lambda = 0$$

$$9\lambda = 6$$

$$\lambda = \frac{2}{3}$$

$$\frac{2}{3} - 8 + 8\lambda = 0$$

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

$$8\lambda = 6$$

$$\lambda = \frac{3}{4}$$

OP1

$$\min_{w, b} \frac{1}{2} \|w\|^2$$

$$\text{s.t. } \operatorname{tr}(w^T \phi(x_n) + b) \geq 1$$

$$\text{solution} = w_1^\top, b_1^\top$$

OP2

$$\min_{w, b} \frac{1}{2} \|w\|^2$$

$$\text{s.t. } \operatorname{tr}(w^T \phi(x_n) + b) \geq 1$$

OP3

$$\min_{w, b} \frac{1}{2} \|w\|^2$$

$$\text{s.t. } \operatorname{tr}\left(\frac{w}{n}^T \phi(x_n) + \frac{b}{n}\right) \geq 1$$

OP3

remove

$$\begin{cases} \min_{w, b} \frac{1}{2} \|w\|^2 \\ \text{s.t. } \operatorname{tr}(w^\top \phi(x_n) + b) \geq 1 \end{cases}$$

OP3 has solution

$$w_2^\top = w_1^\top n$$

$$b_2^\top = b_1^\top n$$

Now, decision hyperplanes are,

$$H1 = \{x | w_1^\top \phi(x) + b_1^\top = 0\}$$

$$H2 = \{x | w_2^\top \phi(x) + b_2^\top = 0\}$$

$$= \{x | nw_1^\top \phi(x) + nb_1^\top = 0\}$$

$$= \{x | w_1^\top \phi(x) + b_1^\top = 0\}$$

$$H2 = H1 \text{ g.r.d}$$

$$w_1 = w_1^\top n$$

$$b_1 = b_1^\top n$$

③ K MCP Kernel multi-class perceptron

1 define $b(\omega) = \sum_{ij} \alpha_{ij} [\phi(x_i, t_i)^T \phi(x_j, t_j) - \phi(x_i, c_j)^T \phi(x_j, t_j)]$

$$\text{if } w^T x = \sum_{n=1}^N \alpha_n t_n x_n^T x = \sum_{n=1}^N \alpha_n t_n K(x_n, x)$$

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

3 for $i = 1 \dots n$

$$4 \quad c_j = \operatorname{argmax}_{t \in T} b(x_i, t) \quad \left. \begin{array}{l} \\ \text{if } h_n = \operatorname{sgn}(b(x_i)) \end{array} \right\}$$

$$5 \quad \text{if } c_j \neq t_i \text{ then } \left. \begin{array}{l} \alpha_{ij} = \alpha_{ij} + 1 \\ \text{and } h_n \end{array} \right\}$$

6

Repeat

Testing: $t^* = \operatorname{argmax}_{t \in T} b(x, t) \quad \left. \begin{array}{l} \\ \text{if } h_n = \operatorname{sgn}(b(x)) \end{array} \right\}$

② MCP (concept of kernel comes from here)

initialize parameters $w = 0$

for $i = 1 \dots N$

$$c_j = \operatorname{argmax}_{t \in T} w^T \phi(x_i, t) \quad \left. \begin{array}{l} \\ \text{if } c_j \neq t_i \text{ then } \end{array} \right\}$$

$$w = w + \phi(x_i, t_i) - \phi(x_i, c_j)$$

w is impinviant and is the weighted average

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

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

① KDP

define $b(\omega) = w^T x = \sum_{ij} \alpha_{ij} K(x_i, x)$

initialize dual parameter $\alpha_{ij} = 0$

for $i = 1 \dots N$

$$\left. \begin{array}{l} h_n = \operatorname{sgn}(b(x)) \\ \text{if } h_n \neq t_i \text{ then } \\ \alpha_{ij} = \alpha_{ij} + 1 \end{array} \right\} \text{Repeat}$$

TEST: $y(\omega) = \operatorname{sign}(b(\omega))$

$$MI(X,Y) = \frac{1}{N} \sum_{i,j} P(X=i, Y=j) \cdot \frac{P(X=i, Y=j)}{P(X=i) \cdot P(Y=j)} = \begin{cases} 0 & \text{when } X, Y \text{ independent} \\ \max & \text{when } X=Y \end{cases}$$

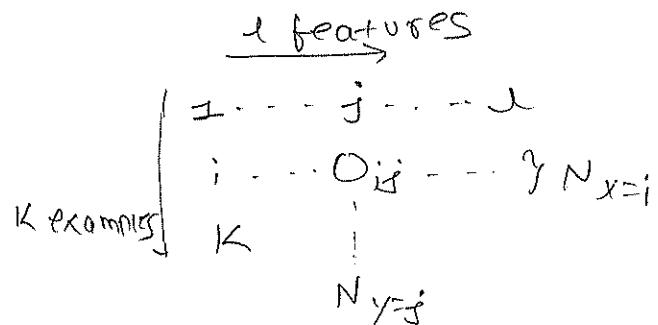
6 mutual information

- works with nominal
- biased towards high purity feature
- may choose redundant feature

② Chi-square

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

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



③ Pearson corr. coeff

$$\rho(X,Y) = \frac{\text{cov}(X,Y)}{\sigma_X \sigma_Y} = \frac{\mathbb{E}[(X-\mu_X)(Y-\mu_Y)]}{\sigma_X \sigma_Y}$$

④ SNR

$$\mu(X, Y) = \frac{\text{binary}}{|H_+ - H_-|} = \frac{|H_+ - H_-|}{\sigma_{\text{sig}} + \sigma_{\text{noise}}}$$

⑤ T-test

$$T(X, Y) = \frac{|H_+ - H_-|}{\sqrt{\frac{\sigma_+^2}{N_+} + \frac{\sigma_-^2}{N_-}}}$$

Filter

VF

Wrapper

- 1 much faster since no need to train the model
- 2 use statistical method of evaluation
- 3 might fail to find best subset
- 4 less prone to overfitting

• computationally expensive

- uses cross-validation
- finds best subset
- more prone to overfitting

① SVM for regression

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

$$\text{s.t. } t_n \leq \bar{t}_n(\bar{w}^T \phi(x_n) + b) + \xi_n$$

$$\bar{t}_n \geq \bar{t}_n(\bar{w}^T \phi(x_n) + b) - \xi_n$$

$$-\xi_n \geq 0 \quad \forall 1 \leq n \leq N$$

② SVM for ranking

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

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

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

MDP Markov Decision process

① Policy

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

(policy expectation)

$$\pi^* = \underset{a}{\operatorname{argmax}} \sum_{s'} T(s, a, s') V(s')$$

(best action policy)

$f = \text{Reward}$

$\gamma t R = \text{discounted reward}$

② Utility

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

(expectation)

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

expectation & max a/s

* Naive Bits

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

$p(y=1)$

$$\begin{array}{ll} P(x_1=1 | y=0) & P(x_1=1 | y=1) \\ P(x_2=1 | y=0) & P(x_2=1 | y=1) \\ P(x_3=1 | y=0) & P(x_3=1 | y=1) \end{array}$$

• # of parameters = $2^{M+1} = 2 \times 3 + 1 = 7$

$$\begin{aligned} \text{• # of parameters if no conditional independence} &= H[2(2^M - 1)] = H[2(8-1)] \\ &= H[2(2^3 - 1)] = H[14] = 15 \end{aligned}$$

Example

category table

Train	DOC	WORD w_i	CLASS	TOTAL
1 chinese 1 2 Beijing 1 3 shanghai 1 4 macao 1 5 Tokyo 1 6 Japan 1	3 documents for C_1 1 document for C_2	1 2 3 4 5 6	chinese Beijing Chinese chinese Chinese Shanghai chinese Macao Tokyo Japan Chinese	$n_1 = 8$ $n_2 = 3$ 11 words 6 unique
Test		5		
Total			$ D = D_1 + D_2 = 3 + 1 = 4$	

- ① Vocabulary, $V = \{ \text{chinese, Beijing, Shanghai, Macao, Tokyo, Japan} \}$
 $|V| = 6$

②

category $C_1 = c$

$$\text{prior } p(c_1) = \frac{|D_1|}{|D|} = \frac{3}{4}$$

③ # of words in class C_1 , $n_1 = 8$

$$P(w_1|c_1) = P(\text{Chinese}|c) = \frac{s+1}{8+6} = \frac{6}{14} = \frac{3}{7}$$

$$P(w_2|c_1) = P(\text{Beijing}|c) = \frac{1+1}{8} = \frac{2}{8}$$

$$P(w_3|c_1) = P(\text{Shanghai}|c) = \frac{2+1}{8+6} = \frac{3}{14}$$

$$P(w_4|c_1) = P(\text{Macao}|c) = \frac{0+1}{8+6} = \frac{1}{14}$$

$$P(w_5|c_1) = P(\text{Tokyo}|c) = \frac{1+1}{8+6} = \frac{2}{14}$$

conditional probability

category $C_2 = j (\text{Japan})$

$$\text{prior } p(c_2) = \frac{|D_2|}{|D|} = \frac{1}{4}$$

of words in class C_2 , $n_2 = 3$

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

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

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

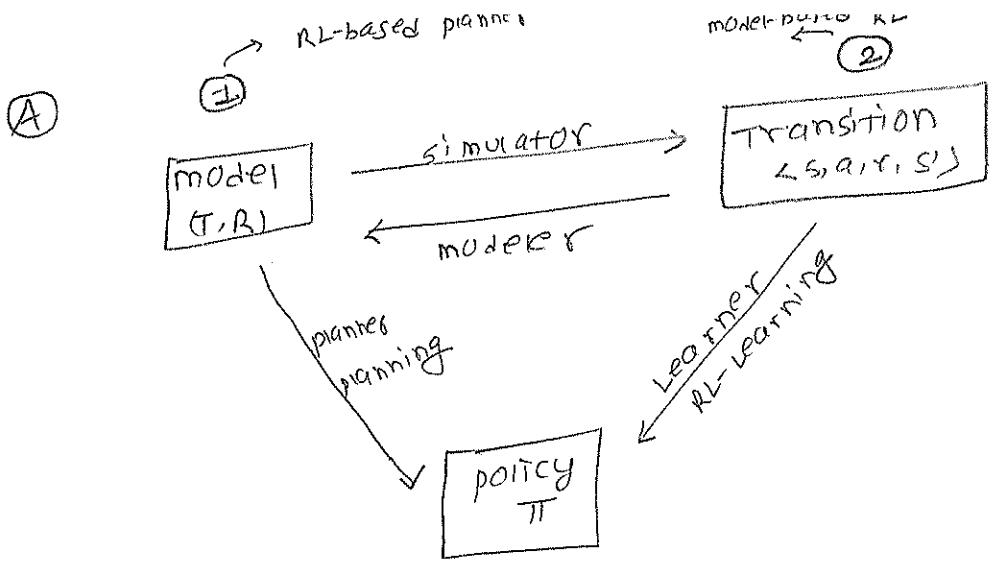
④ choosing a class

$$P(c_1|x) \propto p(c_1) \prod_{i=1}^n P(w_i|c_1) = \frac{3}{4} \cdot \left(\frac{3}{7}\right)^3 \cdot \frac{1}{14} \cdot \frac{1}{14} \approx 0.0003$$

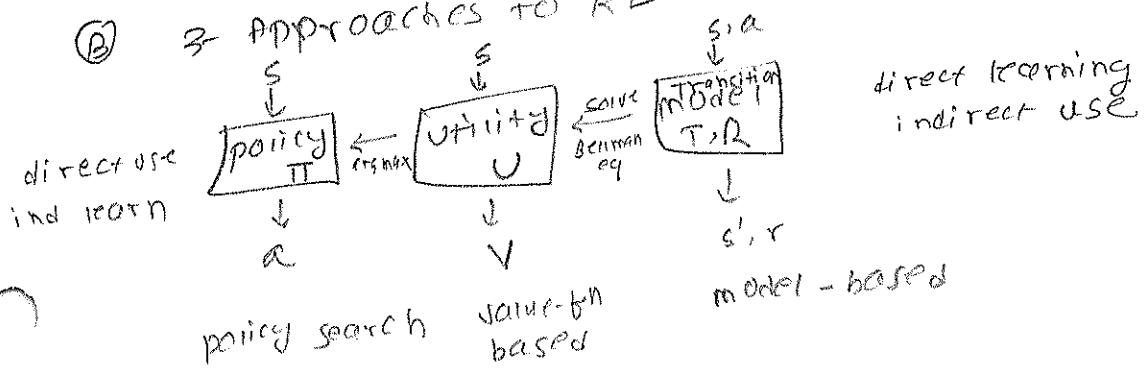
$$P(c_2|x) \propto p(c_2) \prod_{i=1}^n P(w_i|c_2) = \frac{1}{4} \cdot \left(\frac{2}{9}\right)^3 \cdot \frac{2}{9} \cdot \frac{2}{9} \approx 0.0001$$

$$c^* = \operatorname{argmax}_{c_k} P(c_k) \prod_{j=1}^n P(w_j|c_k) = \operatorname{argmax}_{c_k} \{ 0.0003, 0.0001 \}$$

$$-1 \approx -0.0003 = c_1$$



(B) 3 approaches to RL



(C) Q function

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

(Bellman eqn) Utility is scalar

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

(policy gives an action)

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

quiz \Rightarrow $V(s) = \max_a Q(s, a)$

$\Rightarrow \pi(s) = \arg \max_a Q(s, a)$

17 NOV

① mutual information

$$MI(X, Y) = \sum_x \sum_y p(x, y) \ln \left(\frac{p(x, y)}{p(x)p(y)} \right)$$

= 0 when X, Y indp
= max when $X = Y$

$$= KL\left[p(x, y) \parallel p(x)p(y) \right]$$

$b^{ad} \rightarrow$ biased towards high arity features

$b^{ad} \rightarrow$ may choose redundant features

$Wd \rightarrow$ works well with nominal features & labels

①

3 parametric approaches \rightarrow

① discriminant function

{ Fisher's linear disc.

perceptron

SVM

inference and decision
are combined as
single learning
problem

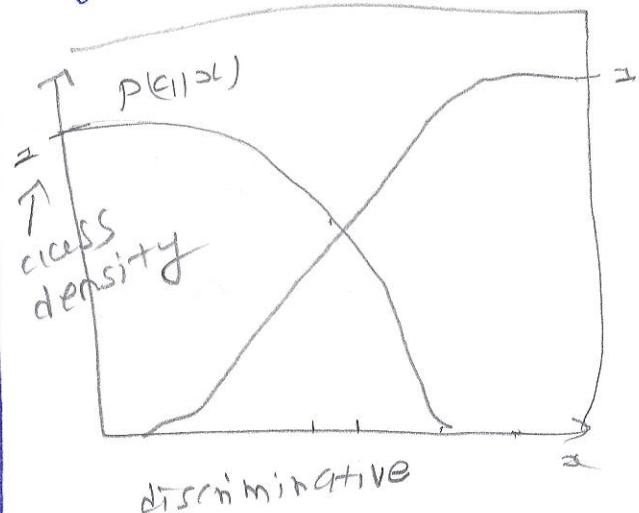
② probabilistic discriminative models

logistic regression
conditional random field

int'l dec \rightarrow separate

others data need to
compute $P(x|y)$
than $P(y|x)$

can accommodate
many overlapping
features



③ Pearson corr coeff

$$S(x,y) = \frac{E[(x-\bar{x})(y-\bar{y})]}{\sigma_x \sigma_y} \quad (\text{population})$$

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

$$-1 \leq S(x,y) \leq 1$$

~~Pearson~~

④ SNR

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

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

μ_+ , σ_+ = mean & std
for +ve class
samples

⑤ T-test

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

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 hr 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.