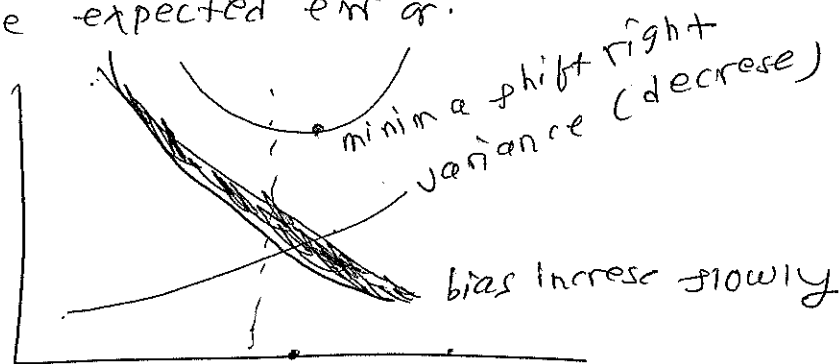


when training data increases, variance is reduced and slightly more complicated model minimizes the expected error.



① show that  $x^T A y$  is a valid kernel if  $A$  is sym. PSD.  
 $\Rightarrow A = Q^T \Lambda Q$  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  $\Lambda = F^T F$

$$\Rightarrow A = Q^T F^T F Q$$

$$\begin{aligned} \Rightarrow x^T A y &= x^T Q^T F^T F Q y \\ &= (F Q x)^T (F Q y) \\ &= \phi(x)^T \phi(y) \end{aligned}$$

$$\text{where } \phi(x) = F Q x$$

Aside:

$$(AB)^T = B^T A^T$$

$$(ABC)^T = C^T B^T A^T$$

$$\downarrow$$

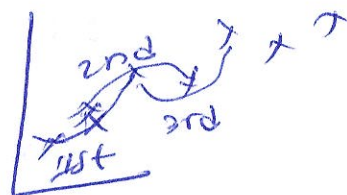
$$((AB)C)^T = C^T (AB)^T = C^T B^T A^T$$

## ⑥ cubic spline smoothing

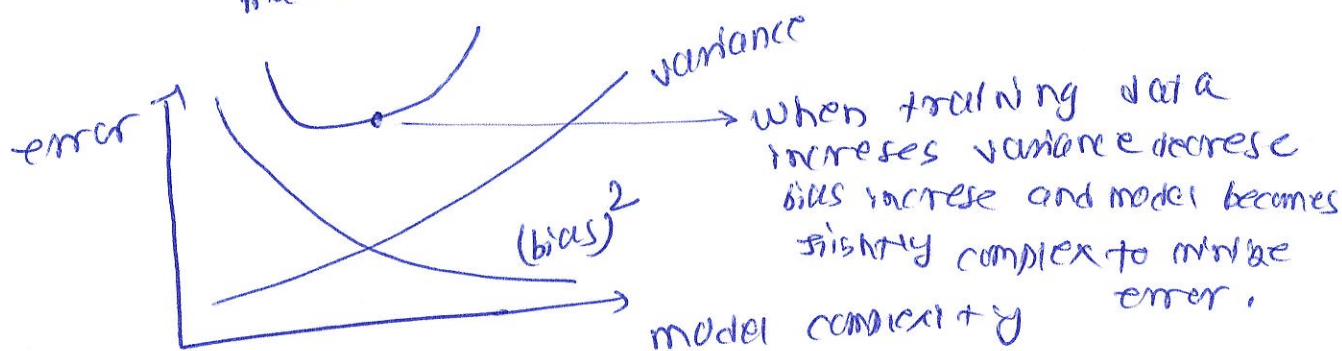
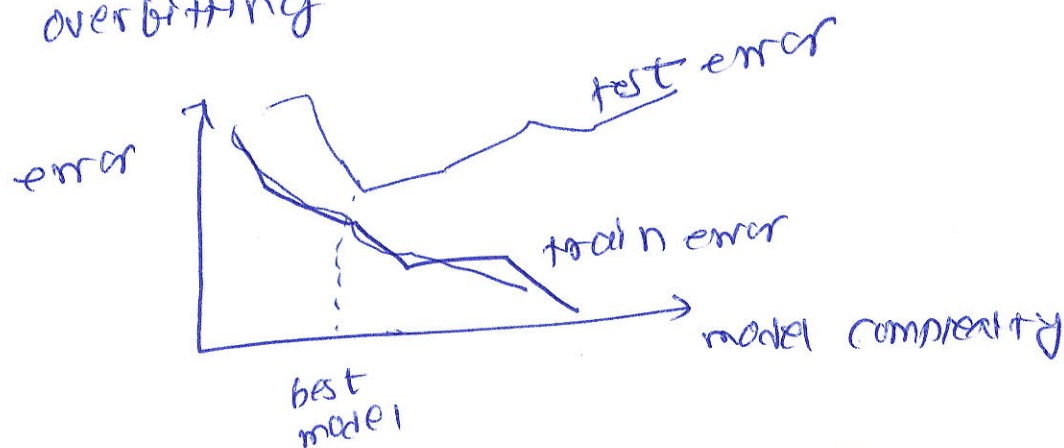
take 3 points  $x, x_{(i)}, x_{(i+1)}$

$$s_i(x) = a_i (x - x_{(i)})^3 + b_i (x - x_{(i)})^2 + c_i (x - x_{(i)}) + d_i$$

$$\forall x \in [x_{(i)}, x_{(i+1)}]$$



## ⑦ overfitting



① Likelihood function for logistic Regression

$$p(h_n | w) = h_n^n \cdot (1-h_n)^{1-n}$$

$$p(w) = \prod_{n=1}^N h_n^{t_n} (1-h_n)^{1-t_n}$$

$$\text{cost } E = \frac{1}{N} \cdot (-\ln p) = -\frac{1}{N} \sum_{n=1}^N [t_n \ln h_n + (1-t_n) \ln (1-h_n)]$$

$$\text{grad} = \nabla_w J = \frac{1}{N} \cdot \sum_{n=1}^N (h_n - t_n) \cdot x_n$$

initial  $z \leftarrow 0, \bar{w} = 0, \tau = 1$   
for  $i = 1 : N$

$h_n = \text{sgn}(w^T x_n)$

if  $h_n \neq t_n$  then

$w = w + t_n x_n$

$\tau = \tau + 1$

$\bar{w} = \bar{w} + w$

$\tau = \tau + 1$

return  $\bar{w} / \tau$

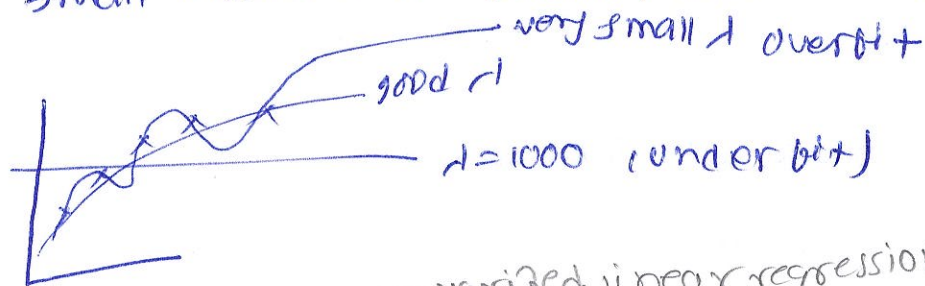
Test:  $\text{sgn}(w^T x)$

② perceptron criterion

$$\text{minimize } E_p(w) = - \sum_{n \in \text{misclassified}} t_n w^T x_n$$

①  $L_2 \rightarrow$  more penalty larger weights  
does not drive weights to zero  
weight vector is NOT sparse.

②  $\lambda \uparrow \Rightarrow$  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  $L_2$ -regularized linear regression

$$③ J = \frac{1}{2N} \sum_{n=1}^N (h_n - t_n)^2 + \frac{\lambda}{2} \|w\|^2$$

$N$ -examples  
in features

$$X = \begin{bmatrix} x_1^1 & x_1^2 & x_1^3 & \dots & x_1^M \\ x_2^1 & x_2^2 & x_2^3 & \dots & x_2^M \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_N^1 & x_N^2 & x_N^3 & \dots & x_N^M \end{bmatrix}$$

$$J(w) = \frac{1}{2N} \sum_{n=1}^N \left( w_0 + \sum_{j=1}^M w_j x_n^j - t_n \right)^2 + \frac{\lambda}{2} \sum_{j=1}^M w_j^2$$

(we omit  $w_0$  in regularization)

regularizing  $w_0$  = shifting origin of target

$\Rightarrow$  some change in all target values  $\Rightarrow$  similar change in estimates

④ Batch GD

stochastic GD

momentum GD  $\hat{w}^{t+1} = \hat{w}^t + \eta \nabla J(\hat{w}^t)$

$$v^{t+1} = \gamma \nabla J(\hat{w}^t)$$

$$\hat{w}^{t+1} = \hat{w}^t - \eta \nabla J(\hat{w}^t)$$

$$v^{t+1}$$

$$\hat{w}^{t+1}$$

$$\gamma \approx 0.9$$

for num-iter:

for num-iter:

for sample in data:

$$w^{t+1} = w^t - \eta \nabla J$$

$$w^{t+1} = w^t - \eta \nabla J$$

learning rate

Nesterov Accelerated Gradient

$$w^{t+1} = w^t - \eta \nabla J(w^t - \gamma \hat{w}^t) - \gamma \hat{w}^t$$

⑤ Gradient descent

$\eta$  = learning rate  
 $\gamma$  = momentum  
 $\theta$  = hyperparameter  
 $\tilde{w}$  = weight vector at time stamp  $t$

$$\rightarrow v^{t+1} = \begin{cases} \eta \nabla J(w^t) & \text{vanilla} \\ \eta \nabla J(w^t) + \gamma v^t & \text{with momentum} \\ \eta \nabla J(w^t - \gamma v^t) + \gamma v^t & \text{Nesterov Accelerated Gradient (NAG)} \end{cases}$$

$$w^{t+1} = w^t - v^{t+1}$$

$\gamma \approx 0.9$

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 of examples and update  $w$  after each

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

$$\nabla J = \frac{1}{N} \sum_{n=1}^N (h_n - t_n) \cdot x_n = \frac{1}{N} \cdot \text{np.sum}((h - t) \cdot x)$$

$$w = w - \eta \text{grad}$$

⑦ GD vs Normal eqn

$$w = w - \eta \text{grad}$$

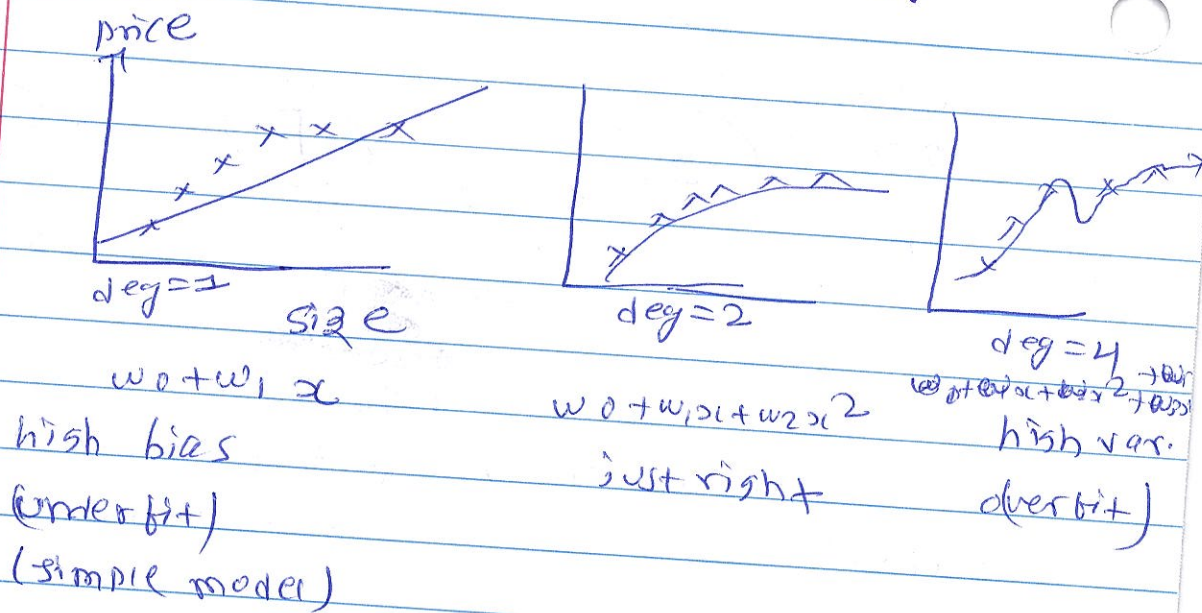
$$w = \underbrace{(X^T X)^{-1}}_{\text{pseudo inverse (pinv)}} X^T t$$

$$w = \underbrace{(X^T X + \lambda N I)^{-1}}_{L2} X^T t$$

for  $L2$  regularized  $I$  is identity matrix of shape  $X^T X$   
 $I[0,0] = 0$  for not to regularize bias term

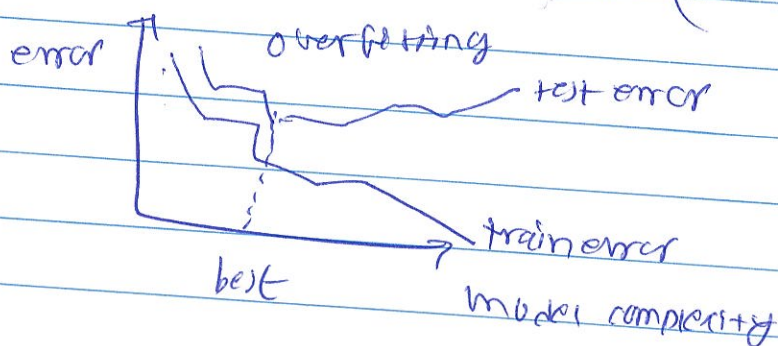
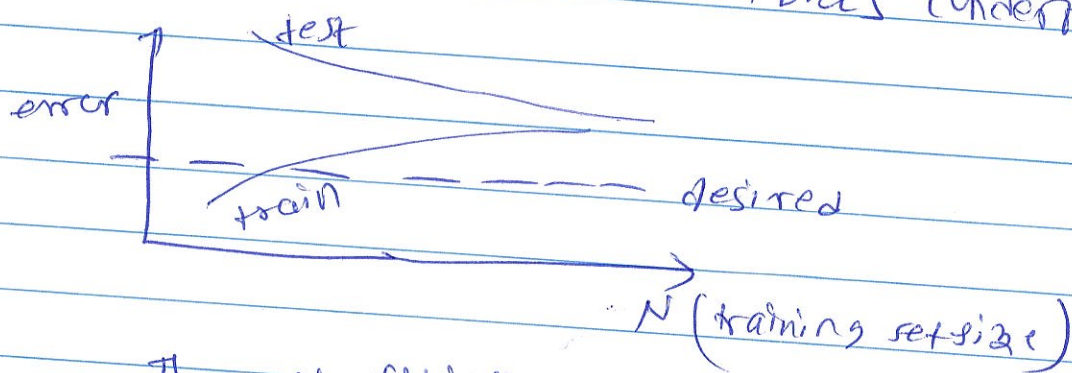


# Bias-variance Trade off

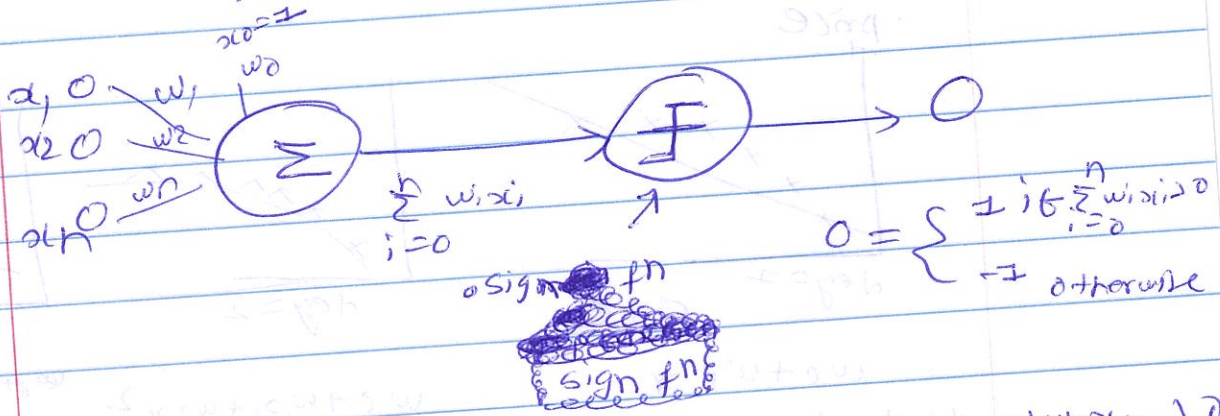


more training eg  
 smaller feature set  $\rightarrow$  bias high  $\sigma^2$  (underfitting)  
 larger feature set  $\rightarrow$  bias high  $\sigma^2$  (")  
 $\rightarrow$  bias high (underfitting)

## Learning curve for high bias (underfit)



# perceptron



output  $O = \begin{cases} \pm 1 & \text{if } \sum_{i=1}^n w_i x_i \geq 0 \\ -1 & \text{otherwise} \end{cases}$

$O = \begin{cases} \pm 1 & \text{if } \vec{w} \cdot \vec{x} \geq 0 \\ -1 & \text{otherwise} \end{cases}$



# ① cost for multivariate linear regression

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

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

$$h = w^T x \quad h = x^T @ w$$

$$h_n = \sum_{j=0}^M w_j x_n^j$$

$$h_1 = w_0 + w_1 x_1^{(1)} + w_2 x_1^{(2)} + w_M x_1^{(M)}$$

first row of sample

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

↑  
bias

$$h = (50, 4) \cdot (4, 1) = (80, 1)$$

$$X = \begin{bmatrix} 1 & 10 & 100 \\ 1 & 20 & 200 \\ 1 & 30 & 300 \\ 1 & 50 & 500 \end{bmatrix}_{50 \times 4}$$

$$t = \begin{bmatrix} 11K \\ 21K \\ 31K \\ 50K \end{bmatrix}_{50 \times 1}$$

2 features ~~age, birth year~~  
age, birth size

$$e = h - t = \begin{bmatrix} h_1 - t_1 \\ h_2 - t_2 \\ \vdots \\ h_N - t_N \end{bmatrix}_{50 \times 1}$$

$$SSE = (h - t) \times x^2$$

$$MSE = \frac{(h - t) \times x^2}{N}$$

$$mse = np.mean(SSE)$$

$$rmse = np.sqrt(mse)$$

$$J = np.sum((h - t) \times x^2) / 2 / N$$

$$J = 0.5 \times np.mean((h - t) \times x^2)$$

$$J = \frac{1}{2} \times mse \quad (\text{single float})$$

$$J = \frac{0.5}{len(t)} (h - t)^T (h - t)$$

$$(h - t)^T (h - t) = \begin{bmatrix} h_1 - t_1 & h_2 - t_2 & \dots \end{bmatrix}_{1 \times 50} \begin{bmatrix} h_1 - t_1 \\ h_2 - t_2 \\ \vdots \\ h_N - t_N \end{bmatrix}_{50 \times 1}$$

$$= \left[ \sum_n (h_n - t_n)^2 \right]_{1 \times 1}$$

## ② gradient of J

$$\nabla_w J = \nabla_w \left[ \frac{1}{2N} \sum_n (h - t)^2 \right] = \begin{bmatrix} \frac{\partial J}{\partial w_0} & \frac{\partial J}{\partial w_1} & \frac{\partial J}{\partial w_2} & \frac{\partial J}{\partial w_M} \end{bmatrix}_{4 \times 1}$$

$$= (h - t) \cdot T @ X^T$$

$$= (1, 50) \cdot (50, 4)$$

$$grad = \nabla_w J = (1, 4)$$

$$grad\_ols = (h - t) \cdot T @ X^T$$

$$\frac{\partial J}{\partial w_j} = \frac{1}{2N} \sum_{n=1}^N (x_n w_j - t_n) \cdot x_n$$

summation variables

$$\nabla_w J = \frac{1}{N} (h - t) \cdot X$$

we use matrix

①  $X = \text{design matrix}$   $\begin{bmatrix} & \end{bmatrix}_{N \times M}$   $N = \text{samples}$   
 $X_1 = \text{biased design matrix}$   $\begin{bmatrix} 1 & \end{bmatrix}_{N, M+1}$   $m = \text{features}$

$t = \begin{bmatrix} \end{bmatrix}_{N \times 1}$  column vector

$w = [w_0 \ w_1 \ w_2 \ w_m]_{1, M+1}$  row vector (2d array)

$w = \text{np.array}(w). \text{reshape}(1, \text{xt.shape}[1]) = (1, M+1)$

$h = X_1 @ w.T = (N, M+1) (M+1, 1) = (N, 1)$  same as  $t$

$e = h - t = N, 1 = 50, 1$

linear regression

$$SSE = \sum_{n=1}^N (h-t)^2$$

$$MSE = \frac{1}{N} \sum_{n=1}^N (h-t)^2$$

$$J = \frac{1}{2N} \sum_{n=1}^N (h-t)^2 = \text{np.sum}((h-t) \times x^2) / 2 / \text{float}(N)$$

$$RMSE = E_{RMS} = \sqrt{MSE}$$

Normal eqn:  $w = (X^T X)^{-1} X^T e$

moore penrose pseudo inverse of  $X$

$\text{np.linalg.pinv}(X)$

# multivariate linear

	bias	x		t
		geography	bed	2K
1		100	2	
1		200	3	3K
				5K
1		500	5	

t //

m features and one bias

$$X = \begin{bmatrix} x_0^{(1)} & x_1^{(1)} & \dots & x_m^{(1)} \\ x_0^{(2)} & x_1^{(2)} & \dots & x_m^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ x_0^{(N)} & x_1^{(N)} & \dots & x_m^{(N)} \end{bmatrix}_{N \times m+1} \quad \begin{bmatrix} t_1 \\ t_2 \\ t_3 \\ \vdots \\ t_n \end{bmatrix}_{N \times 1}$$

There are m features and one bias total (m+1)  
There are N samples for each feature

$$J = \frac{1}{2N} \|XW - t\|^2 = \frac{1}{2N} (XW - t)^T (XW - t)$$

$$\nabla_W J = \frac{1}{N} \frac{d}{dX} (XW - t) \cdot X(XW - t) \quad \frac{d}{dX} (X^T X) = 2X$$

X is a matrix

$$0 = X^T (XW - t)$$

$$X^T X W = X^T t$$

weights

$$W = (X^T X)^{-1} (X^T t)$$

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

derivatives of ~~matrix~~ matrix products

$L_1$  vs  $L_2$  norm  
 $L_2$  more used

$L_2 \rightarrow$  more penalty on large weights, but doesn't drive small weights to zero.

$L_1 \rightarrow$  less penalty for large wt, but leads many weights ~~towards~~ to (or very very close) to zero. leading to weight vector to be sparse.

Bias-Variance Tradeoff

(a) Ridge: 
$$\sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 = RSS + \lambda \sum_{j=1}^p \beta_j^2$$

shrinkage parameter  $\lambda$  which has to be tuned via validation set  $\rightarrow$  regularization strength

Ridge  $\left\{ \begin{array}{l} \lambda \uparrow \Rightarrow \text{var} \downarrow \text{bias} \uparrow \\ \text{meaning: small change in training data} \\ \text{big change in parameter estimates} \\ \text{effect will increase with no. of parameters} \end{array} \right.$

(b) LASSO  $\beta_j^2$  vs  $|\beta_j|$

Ridge shrinks wj but do not make 0  
 but LASSO makes them zero, makes less no of wts.



# Logistic Regression → classification (binary)

linear regression  $h = w^T x$

logistic regression  $h = \sigma(w^T x) = \frac{1}{1 + e^{-w^T x}} = h(x) = y(x)$

likelihood function  $p(t|w) = \prod_n h_n^{t_n} (1 - h_n)^{1 - t_n}$

→ log likelihood,  $E$  or  $J = -\ln p(t|w)$

$$= -\sum_{n=1}^N [t_n \ln h_n + (1 - t_n) \ln (1 - h_n)]$$

∴ cost or Error or loss function

$$E_D = -\sum_n [t_n \ln h_n + (1 - t_n) \ln (1 - h_n)]$$

where  $t_n \in \{0, 1\}$  NOT  $\{-1, 1\}$

add regularizer,

$$E_w = \frac{\lambda}{2} w^T w$$

→ then  $\lambda^2$ -regularized logistic regression cost function

$$E = -\frac{1}{N} \sum_n [t_n \ln h_n + (1 - t_n) \ln (1 - h_n)] + \frac{\lambda}{2} w^T w$$

(cost)  $\lambda$  is regularized

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

$$h_n = h(x_n) = \sigma(x_n) = \frac{1}{1 + e^{-w^T x_n}} \text{ is sigmoid fn}$$

# Softmax Regression (classification / multiclass)

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

$$X = [x_1, x_2, \dots, x_n]$$

$$t_1, t_2, \dots, t_n \in \{1, 2, \dots, K\}$$

$$w_k = [w_{k0}, w_{k1}, w_{k2}, \dots]^T$$

one weight vector per class,

$$p(c_k | x) = \frac{e^{w_k^T x}}{\sum_j e^{w_j^T x}}$$

prob that  
nth example  $x_n$   
has class  $t_n$

$$p(t_n | x_n) = \frac{e^{w_{t_n}^T x_n}}{\sum_j e^{w_j^T x_n}}$$

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,

$$E_D(w) = -\frac{1}{N} \ln L(w)$$

$$= -\frac{1}{N} \ln \prod_{n=1}^N p(t_n | x_n)$$

$$= -\frac{1}{N} \sum_{n=1}^N \ln p(t_n | x_n)$$

$$= -\frac{1}{N} \sum_{n=1}^N \ln \frac{e^{w_{t_n}^T x_n}}{\sum_j e^{w_j^T x_n}}$$

$w_{t_n}$  is weight vector for  
nth example

$w_k$  is weight vector for  
all the examples belonging  
to class  $k$ .

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

## MLE

sample:  $x_1, x_2, \dots, x_N$

prob of  $x$  is  
 $b(x, w)$

joint prob of  $x_1, x_2, \dots, x_N$  is called

likelihood  $L(w)$

$$L(w) = p(x_1, x_2, \dots, x_N | w) = b(x_1, w) \cdot b(x_2, w) \cdot \dots \cdot b(x_N, w)$$

$$L(w) = \prod_{n=1}^N b(x_n, w)$$

# MAP solution for LR

$$p(t|w) = \prod_{n=1}^N h_n^{t_n} (1-h_n)^{1-t_n}$$

$$p(w) = \left(\frac{2}{\pi}\right)^{\frac{m+1}{2}} e^{-\frac{d}{2} w^T w} \propto e^{-\frac{d}{2} w^T w}$$

$$p(w|t) = \frac{p(t|w) p(w)}{p(t)} \propto p(t|w) p(w)$$

$$\text{now, } \hat{w}_{\text{MAP}} = \underset{w}{\operatorname{argmax}} p(w|t)$$

$$= \underset{w}{\operatorname{argmax}} p(t|w) p(w)$$

$$= \underset{w}{\operatorname{argmin}} -\ln p(t|w) -\ln p(w)$$

$$= \underset{w}{\operatorname{argmin}} -\ln \prod_n h_n^{t_n} (1-h_n)^{1-t_n} -\ln e^{-\frac{d}{2} w^T w}$$

$$= \underset{w}{\operatorname{argmin}} -\sum_n \ln(h_n^{t_n} (1-h_n)^{1-t_n}) + \frac{d}{2} w^T w$$

$$= \underset{w}{\operatorname{argmin}} -\sum_n \underbrace{[t_n \ln h_n + (1-t_n) \ln (1-h_n)]}_{E_D(w)} + \underbrace{\frac{d}{2} w^T w}_{E_w(w)}$$

$$\boxed{\hat{w}_{\text{MAP}} = \underset{w}{\operatorname{argmin}} E_D(w) + E_w(w)}$$

ans

# Softmax Regression

$$P(k|x) = \frac{e^{w_k^T x}}{\sum_j e^{w_j^T x}}$$

$$P(t_n|x_n) = \frac{e^{w_{t_n}^T x_n}}{\sum_{j=1}^K e^{w_j^T x_n}}$$

note

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

$$E_D(w) = \frac{1}{N} \sum_n -\ln \prod_n P(t_n|x_n)$$

$$= -\frac{1}{N} \sum_n \ln \frac{e^{w_{t_n}^T x_n}}{\sum_j e^{w_j^T x_n}}$$

$$= -\frac{1}{N} \sum_n \ln \frac{e^{w_{t_n}^T x_n}}{\sum_j e^{w_j^T x_n}}$$

$$E_D(w) = -\frac{1}{N} \sum_n \sum_k \delta_k(t_n) \ln \left( \frac{e^{w_k^T x_n}}{\sum_k e^{w_k^T x_n}} \right)$$

Regularizer

$$+ \frac{\alpha}{2} \sum_{k=1}^K w_k^T w_k$$

$$\frac{\partial E_D}{\partial w_k} = \alpha w_k$$

(no summation)

Let,  $E_n = \sum_k \delta_k(t_n) w_k^T x_n - \sum_k \delta_k(t_n) \ln \left( \sum_k e^{w_k^T x_n} \right)$

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

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

$$\frac{\partial E_n}{\partial w_k} = [\delta_k(t_n) - p(k|x_n)] x_n$$

$$\Rightarrow \frac{\partial E_D(w)}{\partial w_k} = -\frac{1}{N} \sum_n [\delta_k(t_n) - p(k|x_n)] x_n$$

there are K such equations for each classes.

prevent overflow

$$p(k|x) = \frac{e^{w_k^T x}}{\sum_k e^{w_k^T x}}$$

$$c = \max_{1 \leq k \leq K} w_k^T x$$

$$p(k|x) = \frac{e^{(w_k^T x - c)}}{\sum_k \exp(w_k^T x - c)}$$



## 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  $w$  by maximizing likelihood  $P(w | \text{data})$

MAP  $\rightarrow$  maximizes the posterior prob  $P(w | \text{data})$

② Logistic Regression maps inputs to discrete outputs  
" " continuous ".

## ③ PCA & feature selection

Similarity : reduce the dimension of data  
↓ Difference : feature selection finds a subset of features  
PCA produces a smaller <sup>new</sup> set

## perceptron

perceptron criterion:  $w^T x_n \geq 0$  for  $t_n = +1$   
 $w^T x_n < 0$  for  $t_n = -1$

want:  $t_n w^T x_n \geq 0$  for all patterns  
minimize:  $-w^T x_n t_n$  for all misclassified patterns  $n$

$$\boxed{E_P(w) = - \sum_{n \in M} w^T x_n t_n}$$

perceptron ↑ 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.

$$\begin{aligned}K(x, y) &= \langle \phi(x), \phi(y) \rangle \\&= \langle \phi(y), \phi(x) \rangle \quad (\because \text{property of inner product}) \\&= K(y, x)\end{aligned}$$

② show  $K(x, z) = \alpha^T A^T A \beta$  is a valid kernel.

Let  $\phi(x) = Ax$ ,

then,

$$\begin{aligned}\langle \phi(x), \phi(z) \rangle &= \phi(x)^T \phi(z) \\&= (Ax)^T (Az) \\&= \alpha^T A^T A \beta \\&= K(x, z)\end{aligned}$$

$\therefore K(x, z)$  is an inner product in some Hilbert space.

③  $AA^T$  is PSD matrix (indirect)

proof Let,  $A \in \mathbb{R}^{m \times n}$

$\lambda$  = eigenvalue of  $AA^T$

$q$  = eigenvector of  $\lambda$

$$(AA^T)q = \lambda q \quad (\text{premultiply by } q^T)$$

$$\Rightarrow q^T AA^T q = q^T \lambda q$$

$$\therefore AA^T = A^T A$$

$$\begin{aligned} \Rightarrow \lambda &= \frac{q^T AA^T q}{q^T q} = \frac{q^T A^T A q}{q^T q} \\ &= \frac{z^T z}{q^T q} \quad \text{where } z = A^T q \end{aligned}$$

$$\text{here } z^T z \geq 0$$

$$q^T q \geq 0$$

$\therefore \lambda \geq 0$  so  $AA^T$  is PSD.



## MLE vs MAP

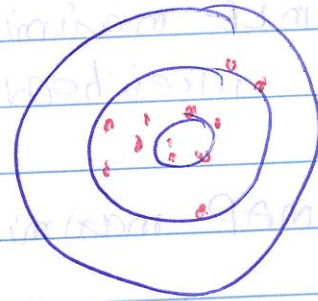
- MLE maximizes the  $\ln p(D|w)$   
likelihood of model  $w$  w.r.t. data  $D$
  - MAP maximizes the  $\ln p(w|D)$   
likelihood of data  $D$  w.r.t. model  $w$
- moreover, MAP additionally uses priors.

QAM 2v - FJM

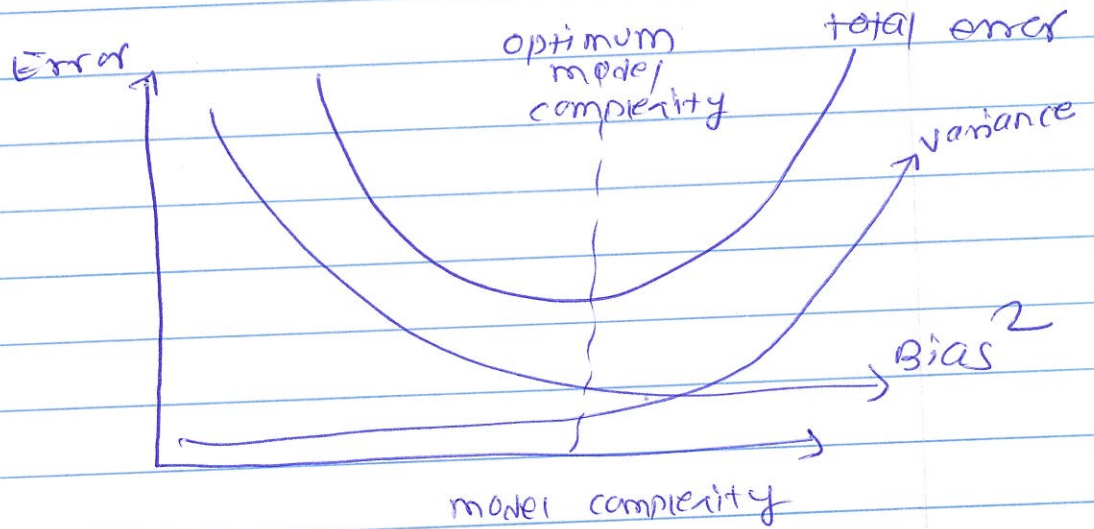
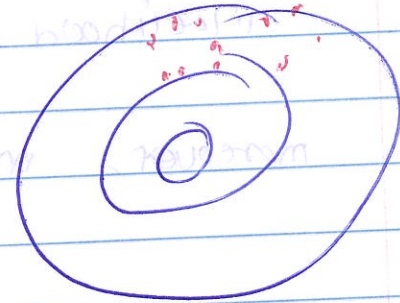
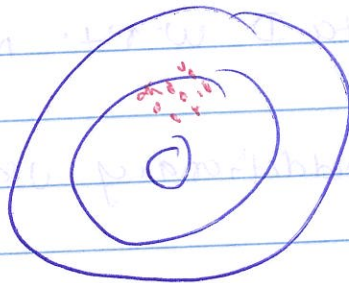
Low variance

High variance

Low Bias



High Bias



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

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

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

$$K(x, z) = (x^T z)^2 = (x_1 z_1 + x_2 z_2)^2 \quad x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

$$= (x_1^2 z_1^2 + 2x_1 z_1 x_2 z_2 + x_2^2 z_2^2) \quad z = \begin{pmatrix} z_1 \\ z_2 \end{pmatrix}$$

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

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

$$\phi(x) = \begin{pmatrix} x_1^2 \\ \sqrt{2} x_1 x_2 \\ x_2^2 \end{pmatrix}$$

$$x^T z = \begin{pmatrix} x_1 & x_2 \end{pmatrix} \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} = x_1 z_1 + x_2 z_2$$

this gives,  $\phi(x) \phi(z) = \begin{pmatrix} x_1^2 z_1^2 + 2x_1 x_2 z_1 z_2 + x_2^2 z_2^2 \end{pmatrix}$

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

\* A necessary and sufficient condition for a kernel function to be "valid" is that the gram matrix be positive and semi-definite for all choices of  $\{\vec{x}_m\}$ .

A gram matrix of  $\vec{x}$  is  $x^T x$ .

The linear vector  $\vec{x}$  is projected into a quadratic surface.

If all the points in this surface are non-zero then our kernel is valid.

\* Fisher's discriminant cost

$$J(w) = \text{Tr}\{W S W^T\}^{-1} \cdot (W S_B W^T)$$

— x — x — x — x — x TUE OCT 31

least square perceptron

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

DOES NOT WORK

$$\min E = 0$$

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

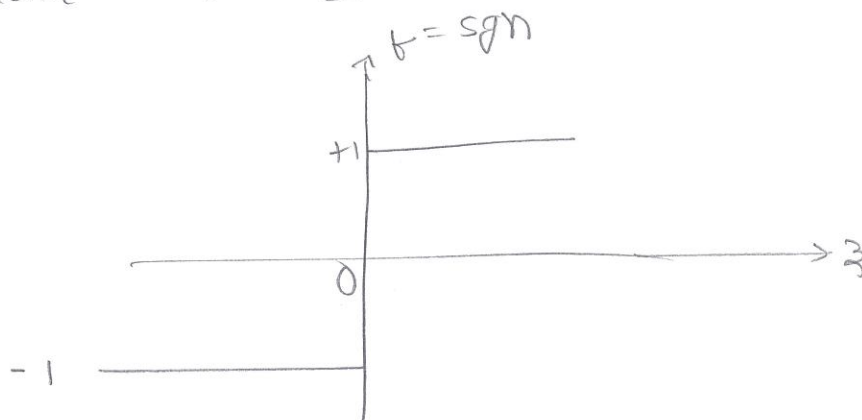
$$0 \rightarrow \frac{1}{2} \cdot \frac{N-2}{2} \rightarrow \frac{N}{2}$$

$N+1$  values

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

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

$$\text{cost} = \left[ 0 \quad \frac{1}{2} \quad \dots \quad \frac{N-2}{2} \quad \frac{N}{2} \right] \text{ discrete set.}$$





# Distance metrics

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

② Hamming  $d(x, y) = \# \text{ of different values in paired length strings}$

③ Mahalanobis distance  $d(x, y) = \sqrt{(x - y)^T S^{-1} (x - y)}$   
 $\uparrow$   
 $S$  is sample cov. matrix

$S = I \rightarrow \text{Euclidean}$

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

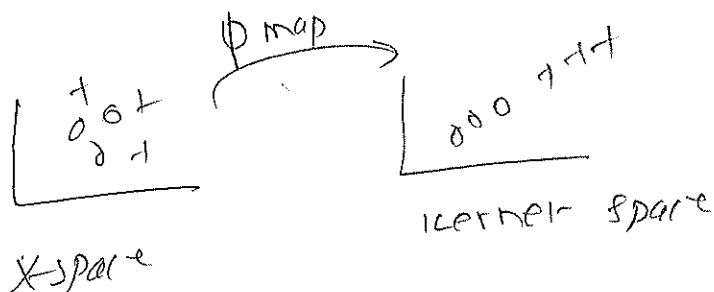
④ cosine similarity

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

⑤ Levenshtein distance (edit distance)

min # of basic operations (del, insert, juxtapose) betn two strings

$x = \text{'attens'}$   $y = \text{'hints'}$   $d(x, y) = 4$



~~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

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

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

i.e.  $1 - \xi_n - t_n w^T \phi(x_n) - b t_n \leq 0$  — ① (convex fn)  
 $-\xi_n \leq 0$  — ② (convex fn constraint)  
 $\forall n = 1, 2, \dots, N$

primal Lagrangian

$$\begin{aligned} L_p = L(w, b, \xi, \alpha, r) &= \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N \xi_n \\ &+ \sum_{n=1}^N \alpha_n (1 - \xi_n - t_n w^T \phi(x_n) - b t_n) \\ &- \sum_{n=1}^N r_n \xi_n \end{aligned}$$

dual Lagrangian

$$L_d(\alpha, r) = \inf_{w, b, \xi} L_p(w, b, \xi, \alpha, r)$$

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

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

NO summation!  
 $0 = \frac{\partial L_p}{\partial \xi_n} = C - \alpha_n - r_n \Rightarrow C = \alpha_n + r_n \quad \forall n \in \{1, 2, \dots, N\}$

then,

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

$$L_D = \sum_n \alpha_n - \frac{1}{2} \sum_{m,n=1}^N \alpha_m \alpha_n t_m t_n K(\mathbf{x}_m, \mathbf{x}_n)$$

then the optimization problem in dual space is,

$$\max_{\alpha} L_D(\mathbf{r}) = \sum_n \alpha_n - \frac{1}{2} \sum_{m,n} \alpha_m \alpha_n t_m t_n K(\mathbf{x}_m, \mathbf{x}_n)$$

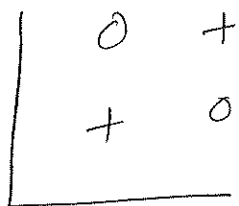
with constraints  $0 \leq \alpha_n \leq C \quad \forall n=1, 2, \dots, N$

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

$$\text{note: } C \sum_n \xi_n - \sum_n \alpha_n \xi_n - \sum_n r_n \xi_n = 0$$

$$\begin{aligned} \text{since, } C \sum_n \xi_n &= \sum_n \alpha_n \xi_n + \sum_n r_n \xi_n \\ &= \sum_n (\alpha_n + r_n) \xi_n \end{aligned}$$

①



- SUM (conv kernel) can achieve zero training error
- Logistic Regr, 3x3 cannot

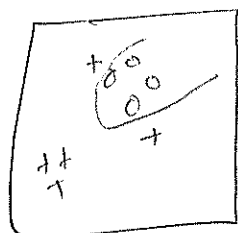
② If examples are iid,  
increase training examples  $\begin{cases} \rightarrow \text{may increase train error} \\ \rightarrow \text{but decrease test error} \end{cases}$

③ SUM effect of  $C$

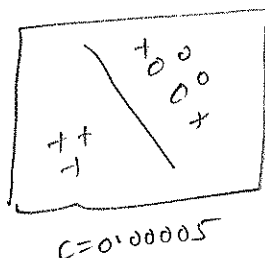
$$\min J(w, b) = \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N \xi_n$$

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

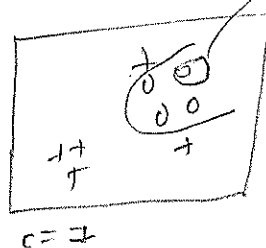
will not change decision bound.



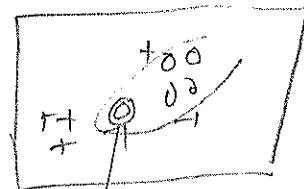
$C = 10,000$



$C = 0.00005$



$C = 1$



adding this change dec. bound, drag it right

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

④ Bias variance Tradeoff

	Bias	var
linear reg	high	low
$d=2$ poly	low	low
$d=10$ poly	low	high

# python

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

$X = \text{np.array}([ [1, 2, 3], [4, 5, 6] ]) = \text{np.arange}(17).reshape(2, 3)$   
 $\text{ones} = \text{np.ones}(X.shape[0]).reshape(-1, 1)$   
 $X_1 = \text{np.append}(\text{ones}, X, \text{axis}=1).astype(\text{np.int})$

$X_1 = \text{np.c_[np.ones(X.shape[0]).reshape(-1, 1), X]$

$X_1 = \text{np.column_stack}([ " ])$

$X_1 = \text{np.c_[np.ones(X.shape[0])[0, np.newaxis], X]$

$X_1 = \text{np.c_[np.ones(X.shape[0]).reshape(-1, 1), X]$

$= \text{np.c_[np.atleast-2d(np.ones(X.shape[0]).T), X]$

$= \text{np.c_[np.expand_dims(np.ones(X.shape[0]), axis=1), X]$

⑤ Given  $\phi(x) = [1, x_1, x_2, x_1x_2]$

and the kernel  $k(x, x')$ .

$$\Rightarrow k(x, x') = 1 + x_1x_1' + x_2x_2' + x_1x_2x_1'x_2'$$

⑥  $L_1$  VS  $L_2$  LOSS

⑨ False:  $L_2$  is more robust to outlier than  $L_1$ .  
gradient of  $L_2$  loss can grow without bounds, but  
gradient of  $L_1$  loss is bounded, hence  
influence of outlier is limited.

(Note:  $L_2$  gives more weight to misclassification  
than  $L_1$ )

⑩  $L_1$  gives sparse solution & used in feature selection.

⑪ Logistic loss is better than  $L_2$  loss in classification task.

⑫ SVM small  $C$ ,

For linearly separable data, small  $C$  can affect  
training accuracy.

A small  $C$  can allow large slack, 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

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

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

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

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

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

compare  $f_i(x) \leq 0$  for  $i = 1, \dots, m$

primal Lagrangian,

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

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  $w, b$ :

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

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

Look LHS

then dual Lagrangian is,

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

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

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

then the optimization problem in dual space is,

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

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

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

Look here

now, KKT conditions are,

① primal constraints  $1 - t_n (\omega^T \phi(x_n) + b) \leq 0$  (convex constraint equation)  
 $t_n (\omega^T \phi(x_n) + b) \geq 1$

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

③ complementary slackness  $\alpha_n \{1 - t_n (\omega^T \phi(x_n) + b)\} = 0$

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

$$1 - t_n \omega^T \phi(x_n) - t_n b = 0$$

these  $\alpha_n$  are called support vectors



pg 6) here  $n=1, 2, \dots, N$

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

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

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

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

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

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

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

$$\left[ b = t_m - \omega \cdot \phi(x_m) \right]$$

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

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

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

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

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

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

$$S \subseteq D$$

$$S = \{n | 1 - t_n \omega^T \phi(x_n) - t_n b = 0\}$$

then, Linear discriminant function is,

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

KNN = memory-based (no model to fit)

① find  $k$  nearest examples  $x_1, x_2, \dots, x_k$  from test set  $x$ .

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

$w_i$  gives distance-weighted KNN

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

② Mahalanobis dist

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

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

KNN usage

- ① satellite image classification
- ② digit recognition

③ cosine similarity

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

④ kernel-based distance weighted NN

→ binary classification  $T \in \{+1, -1\}$

$$y(x) = \operatorname{sign} \left( \sum_{i=1}^N \underbrace{k(x, x_i)}_{\text{kernel}} \cdot t_i \right)$$

## Wrapper method

- ① <sup>greedy</sup> Forward selection
- $F = \text{all features}$
  - $S = \text{subset of features}$
  - start  $S = \{\}$  empty

for each feature  $f$  in  $F - S$   
    find best performing feature  $f$  and add to  $S$ .

Repeat until performance does not increase  
or performance good enough

- ② Recursive Backward Elimination
- $F = \{1, 2, \dots, K\}$  is set of features
  - $S = []$  ranked set of features

Repeat until  $F - S$  is empty

train  $w$  using linear SVM and  $F - S$   
find feature  $f$  with minimum  $|w \cdot f|$   
append  $f$  to  $S$

Return  $S$

### ④ distance-weighted KNN for regression

1. find  $k$  nearest points  $x_1, x_2, \dots, x_k$

$$2. y(x) = \frac{\sum_{i=1}^k w_i t_i}{\sum_{i=1}^k w_i}$$

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

$k=N \rightarrow$  Shepard's method

$$y(x) = \frac{\sum_{i=1}^N K(\|x - x_i\|) t_i}{\sum_{i=1}^N K(\|x - x_i\|)} \quad (\text{kernel-based dist weighted})$$

### ⑤ Regression with KNN

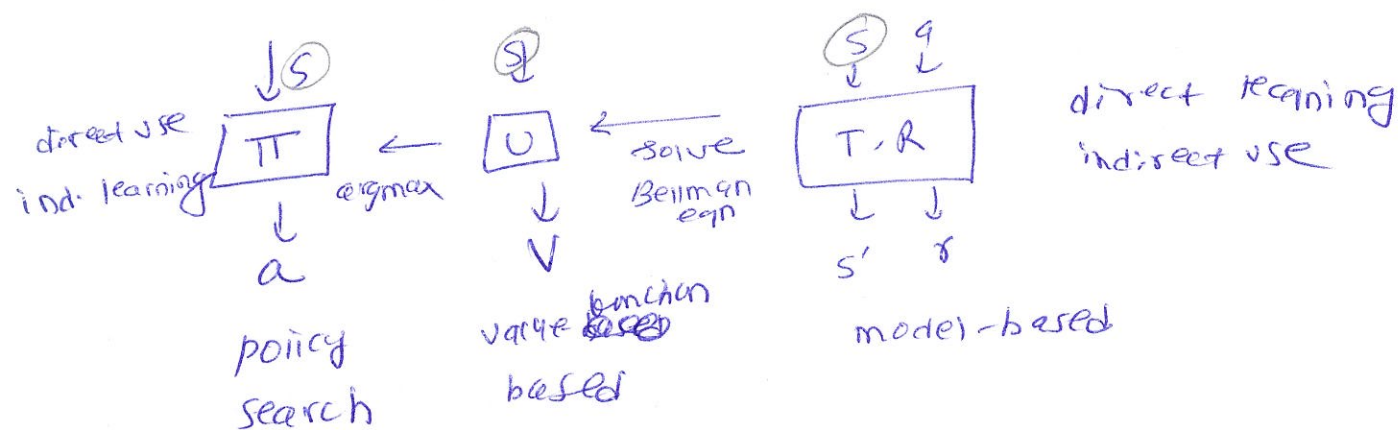
$$y(x) = \frac{1}{k} \sum_{i=1}^k t_i \quad (t_i \text{ are values, not classes})$$

### ⑥ distance-weighted KNN (regression)

original:  $y(x) = \frac{\sum_{i=1}^k w_i t_i}{\sum_{i=1}^k w_i} \quad w_i = \frac{1}{\|x - x_i\|^2}$

kernel based:  $y(x) = \frac{\sum_{i=1}^N K(\|x - x_i\|) t_i}{\sum_{i=1}^N K(\|x - x_i\|)} \quad (t_i \text{ are values not classes})$

### 3 Approaches to RL



### Filter method of Feature Selection

#### ① Mutual Information

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

$\begin{cases} = 0 & \text{when } x, y \text{ indep} \\ \text{max} & \text{when } x = y \end{cases}$

Let there are  $K$  examples with  $L$  features.

②

1	...	$j$	...	$L$
$\vdots$		$O_{ij}$		
$K$		$N_{y=j}$		

$\left. \begin{matrix} \text{ } \\ \text{ } \end{matrix} \right\} N_{x=i}$

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

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

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

classification

③ probabilistic Generative models

{ Naive Bayes

Hidden Markov Models

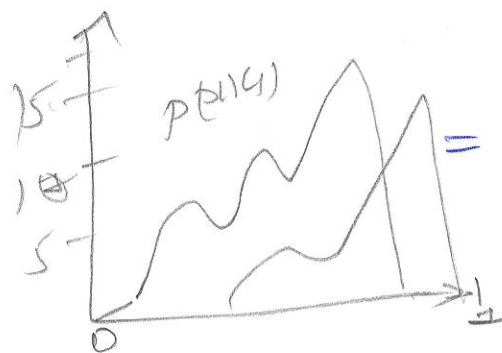
- indices  $\rightarrow$  separate
- can use  $p(x)$  for outlier detection or novelty
- need to model dependencies between features

Naive Bayes  $\rightarrow$  resilient to noise

- text classification with NB  $\rightarrow$  100% HW

- Multiclass probs of class  $c_k$  given ~~test~~ data  $x$  is,

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



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

normalized exponentials

(softmax fn)

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

① SVM for ranking  
optimization problem

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

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

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

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

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

X = np.c\_[np.ones(X.shape[0])[:, np.newaxis].T, X]

② correct = np.sum(y - pred == y - test)  
accuracy = correct / len(y - pred)



⑤ Gradient Descent (GD or BGD) GD with momentum

vanilla GD  $v^{t+1} = \eta \nabla J(w^t)$   
 $w^{t+1} = w^t - v^{t+1}$

$$v^{t+1} = \gamma v^t + \eta \nabla J(w^t)$$

$$w^{t+1} = w^t - v^{t+1}$$

Batch Gradient Descent (Lect 01, p. 24)

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

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

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

$$= w^t - \text{learningRate} (h - t) @ X^T$$

# confusion matrix

ACTUAL cats

	cat	not cat	
predicted cat	TP	FP	$\tilde{P}$
not cat	FN	TN	$\tilde{N}$
	True		
	P	N	

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

	1	0
1	TP	FP
0	FN	TN
	P	N

	actual	
	0	1
0	TN	FN
1	FP	TP
	N	P

$\rightarrow \text{precision} = \frac{TP}{TP+FP}$

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

(hit rate)

$$\frac{\text{predicted true}}{\text{actual true}}$$

F1 score is the harmonic mean of precision and recall,

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

$$= \frac{2 \text{ precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

\* prove that the number of elements in  $X$  and  $Y$  is also a kernel.

$\Rightarrow |X \cap Y|$  is a kernel.

contd (nov 11)

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

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

$$+ \sum_{m \notin S} \alpha_m t_m K(x_m, x)$$

but if  $x_m \notin S$ , then  $\alpha_m = 0$

start  
since  
 $S = \{n \mid \alpha_n > 0\}$   
 $n = 1, 2, \dots, N$   
 $n \notin S \Rightarrow \alpha_n = 0$

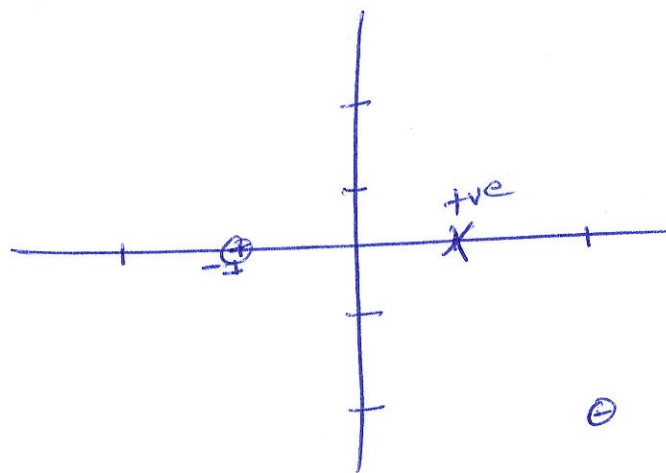
\* package libsvm or SUM Lite

input  $\rightarrow \{x_n, t_n\}_{n=1}^N$   $1 \leq n \leq N$

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

model.txt  $\rightarrow$   $\begin{bmatrix} b \\ \alpha_m t_m \cup x_m \end{bmatrix}$

# perceptron update



$-1, 0$  -ve  
 $2, -2$  -ve

$1, 0 \rightarrow +ve$

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

Augmented data

1	-1	0	-1
1	2	-2	-1
1	1	0	1
$x_{avg}$			$t$

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

note we can  
also use  
 $n \cdot \vec{w} \cdot \vec{x} + b$   
 $n = 0.9$   
 $b = 2$  etc

original unaugmented  $x$

Test point

hypothesis correct or not?

updated weights

$w = w - x = (-1, 1, 0)$

$w = w - x = (-2, -1, 2)$

$w = w + x = (-1, 0, 2)$

$w = w = (-1, 0, 2)$

$w = w = (-1, 0, 2)$

$w = w + x = (0, 1, 2)$

$w = w = (0, 1, 2)$

$w = w = (0, 1, 2)$

$w = w = 0, 1, 2$

Final weight

eg1 - :  $(1, -1, 0) \quad wx = 0 + 0 + 0 < 0$  false

eg2 - :  $(1, 2, -2) \quad wx = -1 + 2 + 0 < 0$  false

eg3 + :  $(1, 1, 0) \quad wx = -2 - 1 + 0 \geq 0$  false

eg1 - :  $(1, -1, 0) \quad wx = -1 + 0 + 0 < 0$  true

eg2 - :  $(1, 2, -2) \quad wx = -1 + 0 + (-4) < 0$  true

eg3 + :  $(1, 1, 0) \quad wx = -1 + 0 + 0 \geq 0$  false

eg1 - :  $(1, -1, 0) \quad wx = 0 + (-1) + 0 < 0$  true

eg2 - :  $(1, 2, -2) \quad wx = 0 + 2 + (-4) < 0$  true

eg3 + :  $(1, 1, 0) \quad wx = 0 + 1 + 0 \geq 0$  true

# ① constrained optimization

(Lagrange multipliers)

$$\text{maximize } 5 - (x_1 - 2)^2 - 2(x_2 - 1)^2$$

$$\text{s.t. } x_1 + 4x_2 = 3$$

sol<sup>n</sup>: If we ignore constraint we get  $x_1 = 2, x_2 = 1$   
then  $x_1 + 4x_2 = 2 + 4(1) = 6$  is too large  
for the constraint.

consider

$$L = L(x_1, x_2, \lambda) = 5 - (x_1 - 2)^2 - 2(x_2 - 1)^2 + \lambda(3 - x_1 - 4x_2)$$

look at:  $\lambda = 0$   $2, 1$

$\lambda = 1$   $3/2, 0$

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

formal sol<sup>n</sup>

$$\frac{\partial L}{\partial x_1} = -2(x_1 - 2) - \lambda = 0$$

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

$$2x_1 = 4 - \lambda$$

$$= 4 - \frac{2}{3}$$

$$= \frac{4 \cdot 3 - 2}{3} = \frac{10}{3}$$

$$x_1 = \frac{5}{3}$$

$$x_2 = \frac{1}{3}$$

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

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

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

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

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

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

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

$$-8 + 8\lambda = 0$$

$$8\lambda = 8$$

$$\lambda = 1/1$$

$$x_1 = 5/3$$

$$x_2 = 1/3$$

$$\lambda = 1$$

OP1

$$\begin{aligned} \min_{w, b} \quad & \frac{1}{2} \|w\|^2 \\ \text{s.t.} \quad & \text{tn}(w^T \phi(x_n) + b) \geq 1 \\ \text{solution} = & w_1^*, b_1^* \end{aligned}$$

OP2

$$\begin{aligned} \min_{w, b} \quad & \frac{1}{2} \|w\|^2 \\ \text{s.t.} \quad & \text{tn}(w^T \phi(x_n) + b) \geq \gamma \end{aligned}$$

$\Downarrow$  OP3

$$\begin{aligned} \min_{w, b} \quad & \frac{1}{2} \left\| \frac{w}{\gamma} \right\|^2 \\ \text{s.t.} \quad & \text{tn}\left(\left(\frac{w}{\gamma}\right)^T \phi(x_n) + \frac{b}{\gamma}\right) \geq 1 \end{aligned}$$

OP3 has solution  $\leftarrow$

$w_2^* = w_1^* \gamma$   
 $b_2^* = b_1^* \gamma$

$\Downarrow$  OP3

$$\begin{aligned} \min_{w, b} \quad & \frac{1}{2} \|w'\|^2 \\ \text{s.t.} \quad & \text{tn}(w'^T \phi(x_n) + b') \geq 1 \end{aligned}$$

where  $w' = w_1^* \gamma$   
 $b' = b_1^* \gamma$

Now, decision hyperplanes are,

$$H_1 = \{x \mid w_1^{*T} \phi(x) + b_1^* = 0\}$$

$$H_2 = \{x \mid w_2^{*T} \phi(x) + b_2^* = 0\}$$

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

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

$$H_2 = H_1 \quad \text{q.e.d}$$

### ③ kmcp kernel multi-class perceptron

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

$\text{kernel } w^T x = \sum_n d_n t_n x^T x = \sum_n d_n t_n K(x_n, x)$

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

3 for  $i = 1 \dots n$

4  $c_j = \arg \max_{t \in T} f(x_i, t)$   
 $(h_i = \text{sgn}(f(x_i)))$

5 if  $c_j \neq t_i$  then  $(h_i \neq t_i)$   
 6  $\alpha_{ij} = \alpha_{ij} + 1$   $(\alpha_{ij} \neq 0)$

Repeat

Testing:  $t^* = \arg \max_{t \in T} f(x, t)$   $(h(x) = \text{sgn}(f(x)))$

### ② mcp (content of kernel comes from here)

initialize parameters  $w = 0$

for  $i = 1 \dots N$

$c_j = \arg \max_{t \in T} w^T \phi(x_i, t)$

if  $c_j \neq t_i$  then  
 $w = w + \phi(x_i, t_i) - \phi(x_i, c_j)$

$w$  is invariant and is the weighted average

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

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

### ① KP

define  $f(x) = w^T x = \sum_n d_n t_n K(x_n, x)$

initialize dual parameter  $d_n = 0$

for  $i = 1 \dots N$

$h_i = \text{sgn}(f(x_i))$

if  $h_i \neq t_i$  then

$d_i = d_i + 1$

Repeat

TEST:  $y(x) = \text{sgn}(f(x))$

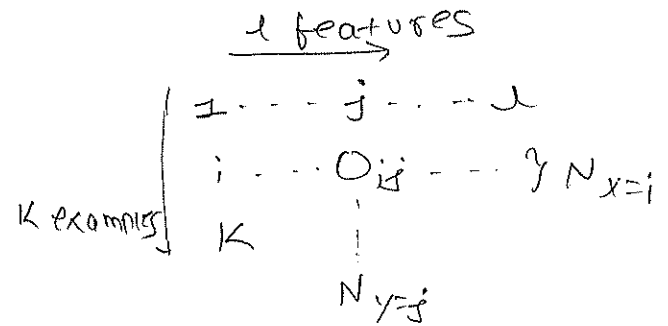
$$MI(X,Y) = \sum_x \sum_y p(x,y) \ln \frac{p(x,y)}{p(x) \cdot p(y)} = \begin{cases} 0 & \text{when } x,y \text{ indep} \\ \text{max} & \text{when } x=y \end{cases}$$

⑥ mutual information  $\begin{cases} \rightarrow \text{works with nominal} \\ \rightarrow \text{biased towards high entropy feature} \\ \rightarrow \text{may choose redundant feature} \end{cases}$

⑦ chi-square

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

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



⑧ pearson corr. coeff

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

⑨ SNR

$$M(x,y) = \frac{|M_+ - M_-|}{\sigma_+ + \sigma_-}$$

$\nearrow$  bin test

⑩ FTest

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



Filter

✓

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 SUM for regression

$$\min J(w, b) = \frac{1}{2} \|w\|^2 + c \sum_{n=1}^N (y_n - \hat{y}_n)^2$$

$$\text{s.t. } t_n \leq t_n(w^T \phi(x_n) + b) + \epsilon + \xi_n$$

$$t_n \geq t_n(w^T \phi(x_n) + b) - \epsilon - \xi_n$$

$$-\xi_n, \xi_n \geq 0 \quad \forall 1 \leq n \leq N$$

② SUM for ranking

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

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

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

# MDP Markov Decision process

① Policy  $\left\{ \begin{array}{l} \pi^* = \underset{\pi}{\operatorname{argmax}} E \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t) \mid \pi \right] \\ \pi^* = \underset{a}{\operatorname{argmax}} \sum_{s'} T(s, a, s') U(s') \end{array} \right.$

(policy expectation)  
(best action policy)

$R$  = Reward  
 $\gamma^t R$  = discounted reward

② Utility  $\left\{ \begin{array}{l} U(s) = E \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t) \mid \pi, s_0 = s \right] \\ U(s) = R(s) + \gamma \max_a \sum_{s'} T(s, a, s') U(s') \end{array} \right.$

(expectation)  
 ← Bellman eqn

## \* Naive Bayes

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

- # of parameters =  $2^{M+1} = 2^{3+1} = 7$
- # of parameters if no conditional independence =  $1 + 2(2^M - 1) = 1 + 2(2^3 - 1) = 1 + 14 = 15$

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

# Example

category label

Train	Doc	words (w <sub>i</sub> )	class	Total
3 documents for c <sub>1</sub>	1	chinese Beijing chinese	c <sub>1</sub>	n <sub>1</sub> = 8
	2	chinese chinese shanghai	c <sub>1</sub>	
	3	chinese macao	c <sub>1</sub>	
1 document for c <sub>2</sub>	4	Tokyo Japan chinese	c <sub>2</sub>	n <sub>2</sub> = 3
Test	5		?	11 words
Total				6 unique

- 1 chinese 1411  
2 Beijing 1  
3 shanghai 1  
4 macao 1  
5 Tokyo 1  
6 Japan 1
- ignore

① vocabulary,  $V = \{\text{chinese, Beijing, shanghai, macao, Tokyo, Japan}\}$   
 $|V| = 6$

②

category  $c_1 = c$

• prior  $P(c_1) = \frac{|D_1|}{|D|} = \frac{3}{4}$

• # of words in class c<sub>1</sub>,  $n_1 = 8$

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

$\rightarrow n_{k1} = n_1 + 1 = 5$   
 $\rightarrow$  Laplace smoothing

$P(w_2|c_1) = P(\text{Beijing}|c) = \frac{1+1}{8+6} = \frac{2}{14} = \frac{1}{7}$

$P(w_3|c_1) = P(\text{shanghai}|c) = \frac{1+1}{8+6} = \frac{2}{14} = \frac{1}{7}$

$P(w_4|c_1) = P(\text{Tokyo}|c) = \frac{0+1}{8+6} = \frac{1}{14}$

$P(w_5|c_1) = P(\text{Japan}|c) = \frac{0+1}{8+6} = \frac{1}{14}$

category  $c_2 = j$  (Japan)

• prior  $P(c_2) = \frac{|D_2|}{|D|} = \frac{1}{4}$

• # of words in class c<sub>2</sub>,  $n_2 = 3$

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

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

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

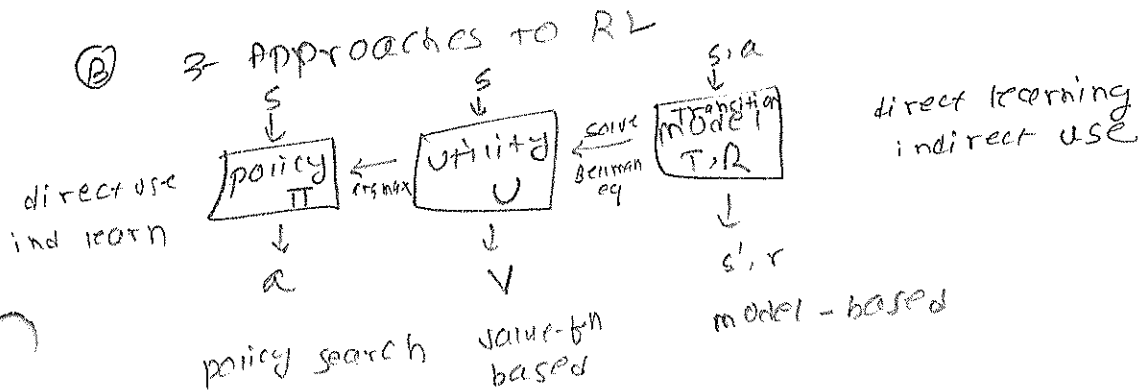
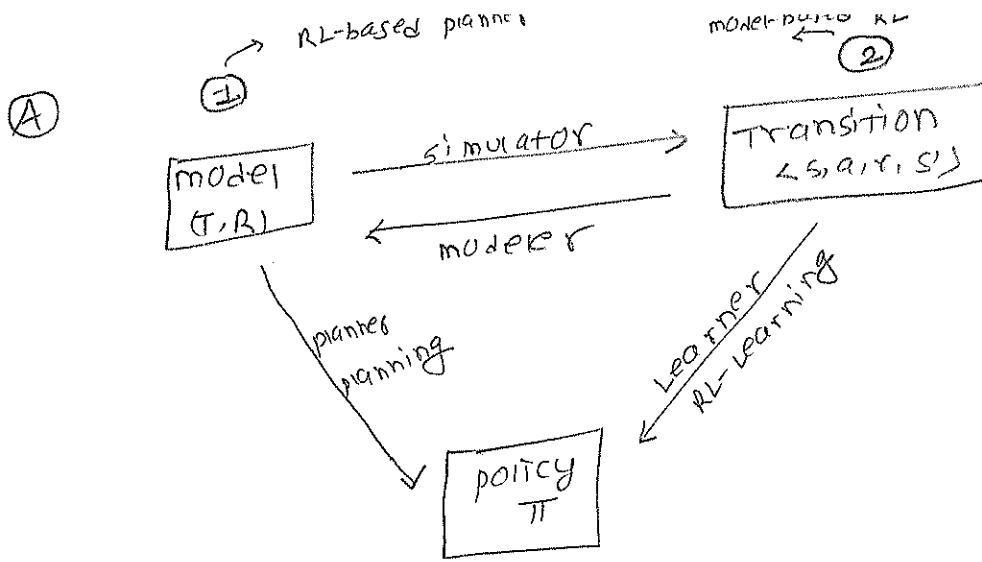
③ choosing a class

$P(c_1|D) \propto P(c_1) \prod_i P(w_i|c_1) = \frac{3}{4} \cdot \left(\frac{3}{7}\right)^3 \cdot \frac{1}{14} \cdot \frac{1}{14} \approx 0.0003$

$P(c_2|D) \propto P(c_2) \prod_i P(w_i|c_2) = \frac{1}{4} \cdot \left(\frac{2}{9}\right)^3 \cdot \frac{2}{9} \cdot \frac{2}{9} \approx 0.0001$

$C^* = \underset{C_k}{\operatorname{argmax}} P(c_k) \prod_{j=1}^n P(w_j|c_k) = \underset{C_k}{\operatorname{argmax}} \{0.0003, 0.0001\}$

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



(C) Q function

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

change  $U$  to  $Q$

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

(policy gives an action)

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

quiz  $\Rightarrow U(s) = \max_a Q(s, a)$

$\Rightarrow \pi(s) = \arg \max_a Q(s, a)$

C

C

C

END

① mutual information

$$\begin{aligned} MI(X, Y) &= \sum_x \sum_y p(x, y) \cdot \ln \left( \frac{p(x, y)}{p(x) p(y)} \right) = 0 \text{ when } X, Y \text{ indep} \\ &= KL[p(x, y) \parallel p(x) p(y)] = \max \text{ when } X = Y \end{aligned}$$

bad  $\rightarrow$  biased towards high arity features

bad  $\rightarrow$  may choose redundant features

bad  $\rightarrow$  works only with nominal features + labels

# \* 3 parametric approaches to

## ① discriminant function

{ Fisher's linear disc.  
perceptron  
SVM

inference and decision  
are combined as  
single learning  
problem

## ② probabilistic discriminative models

Naïve Bayes regression  
conditional random field

• int. dec → separate

• test data need to  
compute  $P(x|c)$

than  $P(c|x)$

• can accommodate  
many overlapping  
features





③ pearson corr coeff

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

(population)

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

~~③~~  $-1 \leq \rho(x, y) \leq 1$

④ SNR

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

binary classes =  $\{+, -\}$   
mean =  $\{\mu_+, \mu_-\}$

examples are binary

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

$\mu_+, \sigma_+ = \text{mean \& std for +ve class samples}$

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 student's 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.