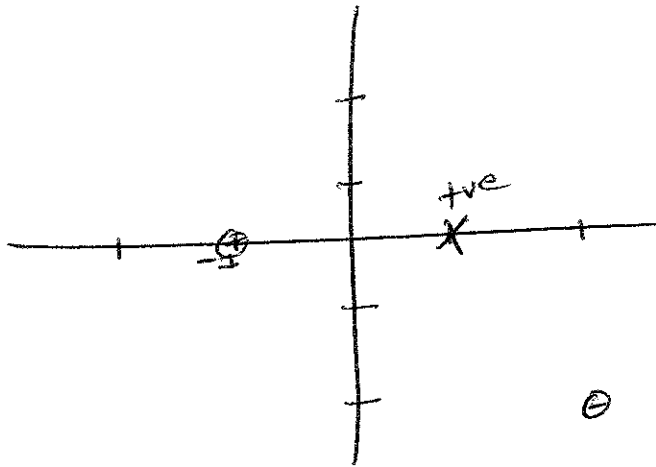


perceptron update



$-1, 0$ -ve
 $2, -2$ -ve
 $1, 0$ → +ve

initial $w = 0, 1, 2$
 let's say new point
 $x = 2, -1.01$
 $x_{avg} = 1, 2, -1.01$
 $w x_{avg} = 0 + 2 - 2.02$
 < 0
 ∴ -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
 $\vec{w} \cdot \vec{x} + b$
 $n = 0.9$
 $b = 1$ etc

Test point

hypothesis
 correct
 misclassified?

updated weights correct example

- eg1 - : $(1, -1, 0)$ $w x = 0 + 0 + 0 < 0$ false
- eg2 - : $(1, 2, -2)$ $w x = -1 + 2 + 0 < 0$ false
- eg3 + : $(1, 1, 0)$ $w x = -2 - 1 + 0 \geq 0$ false
-
- eg1 - : $(1, -1, 0)$ $w x = -1 + 0 + 0 < 0$ true
- eg2 - : $(1, 2, -2)$ $w x = -1 + 0 + (-4) < 0$ true
- eg3 + : $(1, 1, 0)$ $w x = -1 + 0 + 0 \geq 0$ false
-
- eg1 - : $(1, -1, 0)$ $w x = 0 + (-1) + 0 < 0$ true
- eg2 - : $(1, 2, -2)$ $w x = 0 + 2 + (-4) < 0$ true
- eg3 + : $(1, 1, 0)$ $w x = 0 + 1 + 0 \geq 0$ true

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



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	

	<u>actual</u>		
	0	1	
0	TN	FN	\tilde{n}
1	FP	TP	\tilde{p}
	N	P	

$\rightarrow \text{precision} = \frac{TP}{TP + FP}$
 \uparrow
 $\text{recall} = \frac{TP}{P} = \frac{TP}{TP + FN}$

F1 score is the harmonic mean of precision and recall,

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

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

(hit rate)
 $\frac{\text{predicted true}}{\text{actual true}}$

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

$\Rightarrow |K(X, Y)| = |X \cap Y|$ is a kernel.

contd (Nov 11)

$$\begin{aligned} \omega(\phi) &= \sum_{n=1}^N \alpha_n t_n K(x_n, x) \\ &= \sum_{m \in S} \alpha_m t_m K(x_m, x) \end{aligned}$$

start
since \downarrow
 $S = \{n \mid \alpha_n > 0\}$
 $n \notin S \Rightarrow \alpha_n = 0$

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

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

* package libsvm or SUMMIT

input $\rightarrow \{x_n, t_n\} \quad 1 \leq n \leq N$

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

model: $t \mapsto \begin{bmatrix} b \\ \alpha_m t_m \cup x_m \end{bmatrix}$

pg 6) 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.

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

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

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

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

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

$$\therefore b = \frac{1}{t_m} \left(1 - \sum_n \alpha_n t_n \phi(x_n)^T \phi(x_m) \right) = t_m - \sum_n \alpha_n t_n \phi(x_n)^T \phi(x_m)$$

$$b = t_m - \omega \cdot \phi(x_m)$$

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

$$\therefore \frac{1}{t_m} = t_m$$

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

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,

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

where S is the subset of all the examples where lagrange parameter α is non-zero.

$$S \subseteq \mathcal{D}$$

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

then, Linear discriminant function is,

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

then dual Lagrangian is,

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

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

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

then the optimization problem in dual space is,

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

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

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

now, KKT conditions are,

① primal constraints

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

note: we can write $y(x_n) = \omega^T \phi(x_n) + b$

(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

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

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

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

these α_n are called support vectors

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 :

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

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

Look LHS

① constrained optimization

(Lagrange multipliers)

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

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

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

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

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

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

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

✓

formal solⁿ

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

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

$$2x_1 = 4 - \lambda$$

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

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

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

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

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

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

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

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

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

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

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

$$-6 + 9\lambda = 0$$

$$9\lambda = 6$$

$$\lambda = \frac{2}{3}$$

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

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

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

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

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

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



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

let $\phi(x) = Ax$,

then,

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

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

③ AA^T is PSD matrix (indirect)

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

λ = eigenvalue of AA^T

q = eigenvector of λ

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

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

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

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

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

$$q^T q \geq 0$$

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

multivariate Linear

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

t //

m features and one bias

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

$$t = \begin{bmatrix} t_1 \\ t_2 \\ t_3 \\ \vdots \\ t_n \end{bmatrix}_{N \times 1}$$

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

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

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

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

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

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 β_j^2 vs $|\beta_j|$

Ridge shrinks β_j but does not make ϕ
but Lasso makes them zero, makes less no. of wtf.

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

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

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

Batch Gradient Descent (Lect 01, p. 24)

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

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

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

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

④ $X = \text{design matrix } \begin{bmatrix} & \end{bmatrix}_{N \times M}$ $N = \text{samples}$
 $X_1 = \text{biased design matrix } \begin{bmatrix} 1 & \end{bmatrix}_{N, M+1}$ $m = \text{features}$

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

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

$w = \text{np.array}(w). \text{reshape}(1, \text{x1.shape[1]}) = (1, M+1)$

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

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

linear
regression

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

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

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

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

normal eqn: $w = \underbrace{(X^T X)^{-1}}_{\text{moore penrose pseudo inverse of } X} X^T e \ t$

moore penrose pseudo inverse of X

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

① cost for multivariate linear regression

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

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

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

$$t = \begin{bmatrix} 11K \\ 2K \\ 3K \\ 50K \end{bmatrix}_{50 \times 1}$$

2 features ~~age, floor size~~ ~~price~~

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

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

$$h \pm = \omega_0 + \omega_1 x_1^{(1)} + \omega_2 x_1^{(2)} + \dots + \omega_M x_1^{(M)}$$

first row of sample

$$w = \begin{bmatrix} \omega_0 & \omega_1 & \omega_2 & \omega_3 \end{bmatrix}_{1 \times 4}$$

↑
bias

$$h = (50, 4) \cdot (4, 1) = (20, 1)$$

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

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

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

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

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

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

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

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

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

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

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

② gradient of J

$$\nabla_w J = \nabla_w \frac{1}{2N} \sum_{n=1}^N (h_n - t_n)^2 = \begin{bmatrix} \frac{\partial J}{\partial \omega_0} & \frac{\partial J}{\partial \omega_1} & \frac{\partial J}{\partial \omega_2} & \frac{\partial J}{\partial \omega_M} \end{bmatrix}_{4 \times 1}$$

$$= (h - t) \cdot T @ X \pm / N$$

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

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

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

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

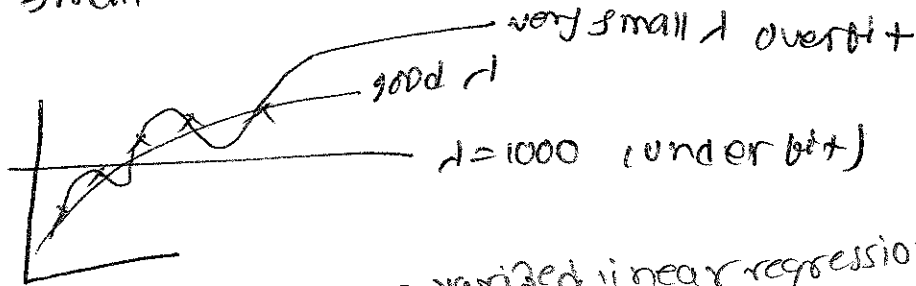
summation variables

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

we use matrix

① $L_2 \rightarrow$ more penalize 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

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

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 same change in all target values \Rightarrow similar change in estimates

④ Batch GD

stochastic GD

momentum & β in \tilde{w}

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

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

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

for num-iter:

for num-iter:

for sample in data:

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

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

$$\eta \approx 0.9$$

learning rate

Nesterov Accelerated
gradient

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

③ Gradient descent

η = learning rate
 v = momentum
 α = hyperparameter
 w = weight vector at time stamp

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

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

Accelerated gradient NAG

② 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

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

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

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

⑤ GD vs normal eqn

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

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

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

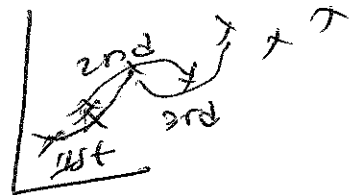
for regularized I is identity matrix of shape $X^T X$
 $I[0][0] = 0$ for not to regularize bias term

⑥ cubic spline smoothing

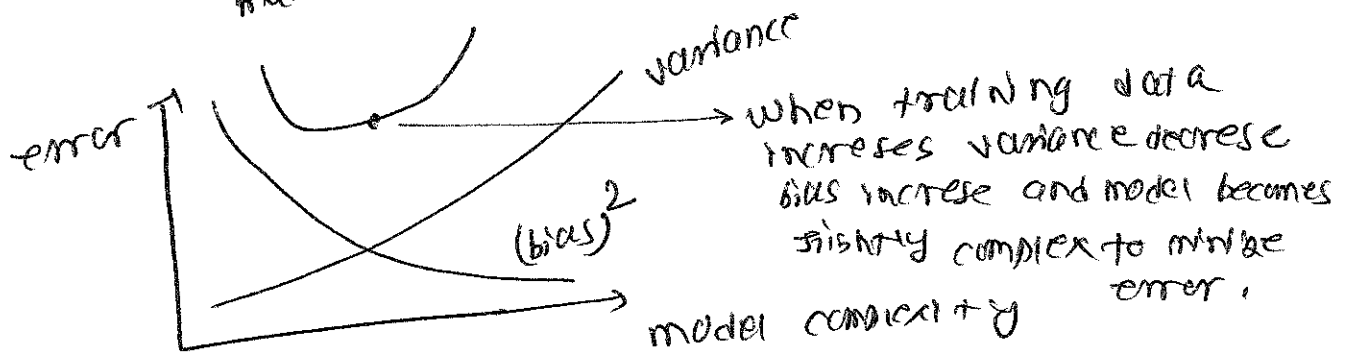
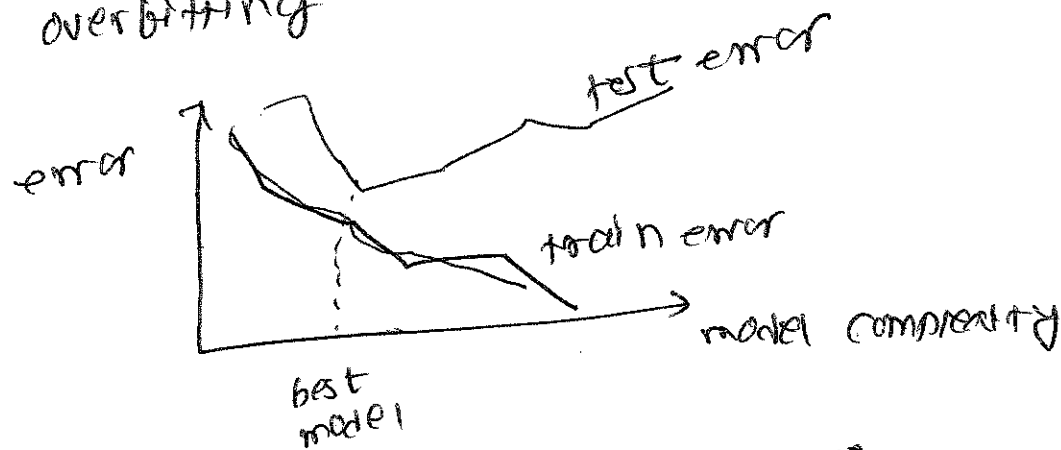
take 3 points $x, x(i), x(i+1)$

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

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



⑦ overfitting



① Likelihood function for logistic Regression

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

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

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

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

② perceptron criterion

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

initialize $w_0, \bar{w} = 0, \tau = 1$
 for $i = 1 : N$
 $h_n = \text{sgn}(w^T x_n)$
 if $h_n \neq t_n + \text{ben}$
 $w = w + t_n x_n$
 $\tau = \tau + 1$
 $\bar{w} = \bar{w} + w$
 $\tau = \tau + 1$
 return \bar{w} / τ
 test: $\text{sgn}(w^T x)$

~~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

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

dual Lagrangian

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

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

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

NO summation!

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

Then,

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

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

Then the optimization problem in dual space is,

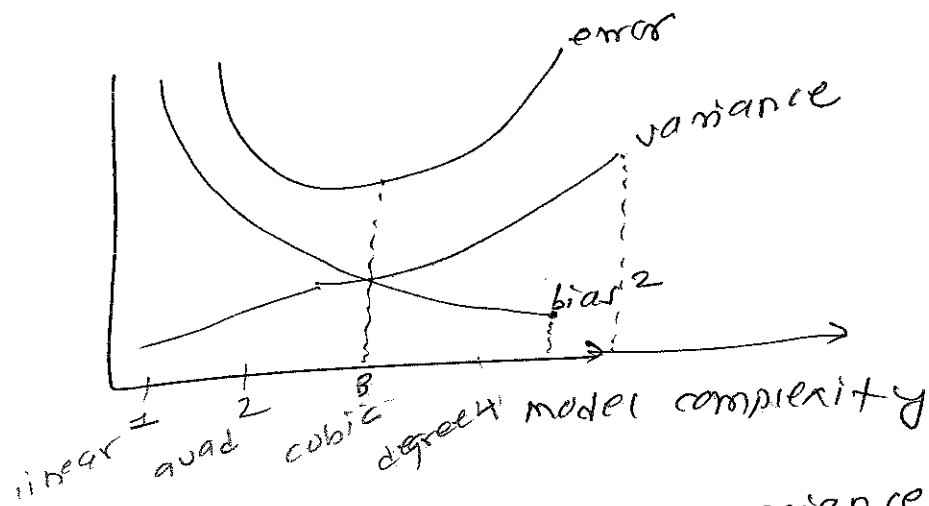
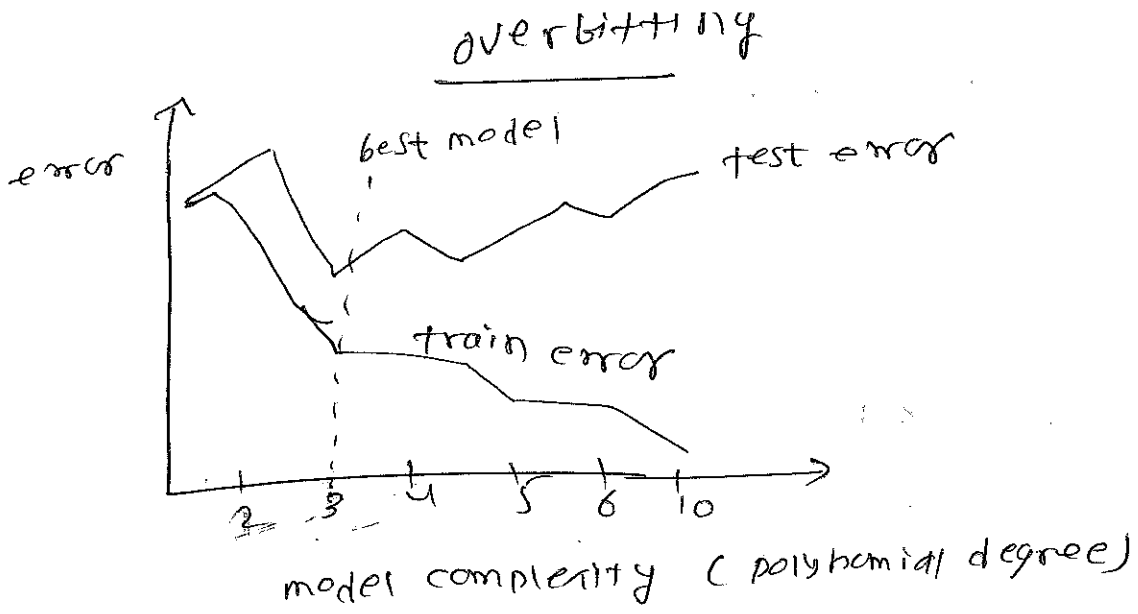
$$\max_{\alpha} L_D(K), = \sum_n \alpha_n - \frac{1}{2} \sum_{m, n} \alpha_m \alpha_n t_m t_n K(\phi_m, \phi_n)$$

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

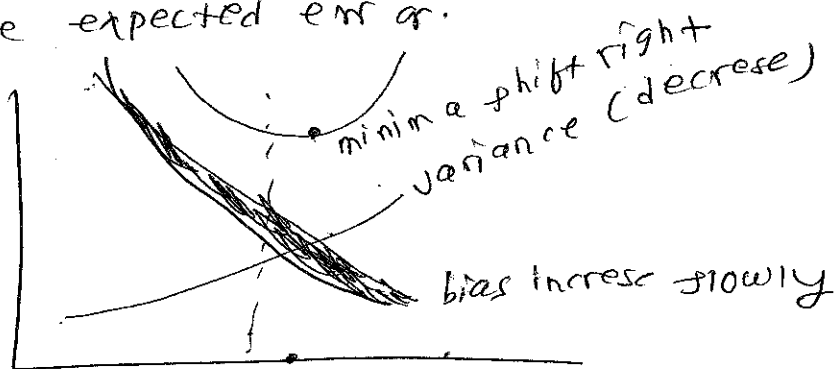
$$\sum_{n=1}^N \alpha_n t_n = 0$$

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

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



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

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

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

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

Aside:

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

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

$$\downarrow$$

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

KNN = memory-based (no model to fit)

① find k nearest examples t_1, t_2, \dots, t_k from test set T

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

w_i gives distance-weighted KNN
 $w_i = \frac{1}{\|x - t_i\|^2}$

② Mahalanobis dist

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

sample covariance matrix

if $S = I$ = Euclidean distance

$$S = \text{diag}(\sigma_1^{-2}, \sigma_2^{-2}, \dots, \sigma_k^{-2})$$

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

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

Wrapper method

- ① ~~greedy~~ 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 ω using linear SVM and $F - S$

 find feature f with minimum $|w_f|$

 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

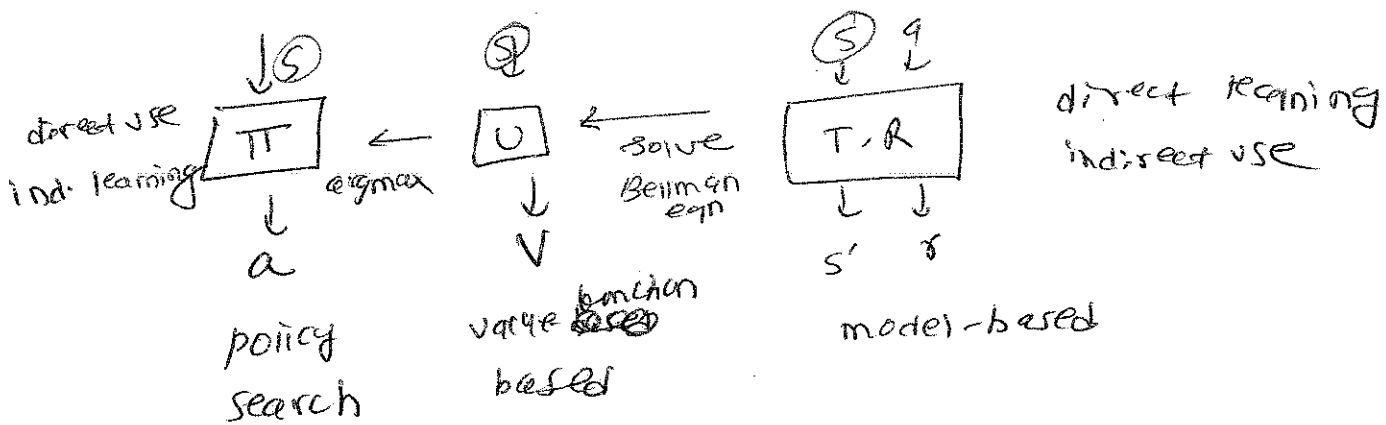
$$y(x) = \frac{1}{k} \sum_{i=1}^k t_i \quad (t_i \text{ are values, not classes})$$

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

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

EQ when x, y indep
max when $x=y$

Let there are K examples with L features.

②

1	...	j	...	L
⋮		O_{ij}		$\} N_{x=i}$
K		$N_{y=j}$		

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

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

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

③ Pearson corr coeff

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

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

~~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
 $y \in \{+1, -1\}$
 $\mu_+, \sigma_+ = \text{mean \& std for two class samples}$

⑤ T-test

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

* 3 parametric Approaches to

① discriminant function

Fisher's linear-DISC
perception
SVM

inference and decision
are combined as
single learning
problem

② probabilistic discriminative models

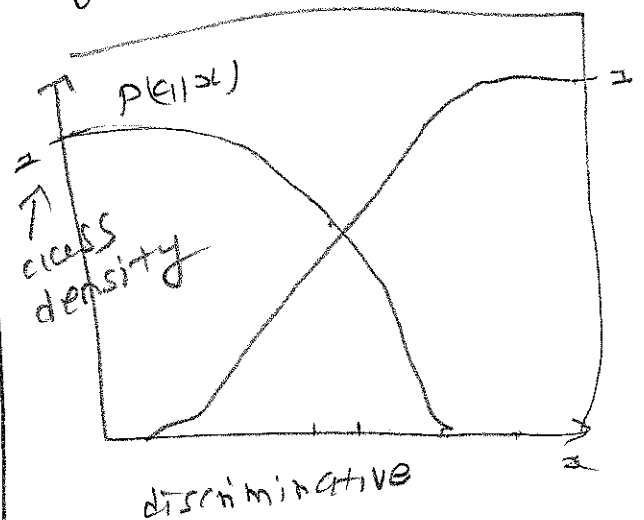
Naive Bayes
Sigmoid regression
conditional random field

inb-dec \rightarrow separate

others data need to
compute $p(x|c)$

than $p(c|x)$

can accommodate
many overlapping
features



classification

② probabilistic Generative models

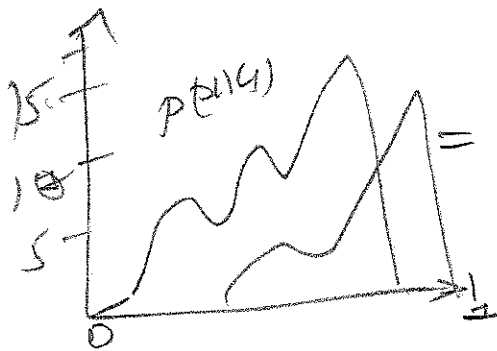
{ Naive Bayes
Hidden Markov Models

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

Naive Bayes \rightarrow resilient to noise

- text classification with NB \rightarrow 100% HW
- Multiple classes
posterior prob 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 $a_{kj}(x) = \ln P(x | c_k) \cdot P(c_k)$

① SVM for ranking
optimization problem

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

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

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

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

Logistic Regression with cross-entropy (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 $p(t|w) = \prod_n h_n^{t_n} (1 - h_n)^{1 - t_n}$

-ve 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

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

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

add regularizer,

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

then L_2 -regularized logistic regression cost function

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

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

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

Softmax Regression

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,

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

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

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

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

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

probability of x is
 $f(x, w)$

joint prob of x_1, x_2, \dots, x_N is called
likelihood $L(w)$

$$L(w) = p(x_1, x_2, \dots, x_N | w) = \prod_{n=1}^N f(x_n, w)$$

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

python

① $X = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}$

append ones column to data X.

$$X = \text{np.array}([[1, 2, 3], [4, 5, 6]]) = \text{np.arange}(17).reshape(2, 3)$$

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

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

$$X1 = \text{np.c_[np.ones}(X.\text{shape}[0]).\text{reshape}(-1, 1), X]$$

$$X1 = \text{np.column_stack}([" "])$$

$$X1 = \text{np.c_[np.ones}(X.\text{shape}[0]).\text{reshape}(-1, 1), X]$$

$$X1 = \text{np.c_[np.ones}(X.\text{shape}[0]).\text{reshape}(-1, 1), X]$$

$$= \text{np.c_[np.atleast_2d}(\text{np.ones}(X.\text{shape}[0]).\text{reshape}(-1, 1)), X]$$

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

①

$$\begin{vmatrix} 0 & + \\ + & 0 \end{vmatrix}$$

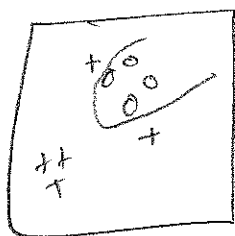
- SUM quad kernel can achieve zero training error
- Logistic Regr, SVM cannot

② If examples are iid, increase training examples \rightarrow may increase train error but decrease test error

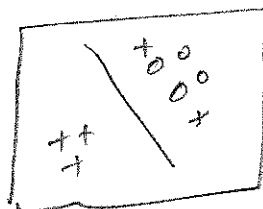
③ SUM effect of C

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

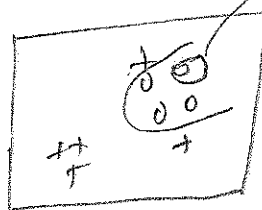
s.t. $w^T \phi(x_n) + b \geq 1 - \xi_n \quad \forall n \in \{0, N\}$



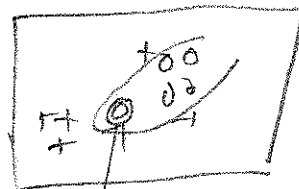
$C = 10,000$



$C = 0.00005$



$C = 1$



adding this change dec. boud. drastically

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

④ Bias variance Tradeoff

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

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

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

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

⑥ $L1$ VS $L2$ LOSS

⑨ False: $L2$ is more robust to outlier than $L1$.
gradient of $L2$ loss can grow without bounds, but,
gradient of $L1$ loss is bounded, hence
influence of outlier is limited.

(Note: $L2$ gives more weight to misclassification
than $L1$)

⑩ $L1$ gives sparse solution & used in feature selection.

⑪ Logistic loss is better than $L2$ 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.

C

C

C

Distance metrics

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

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

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

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

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

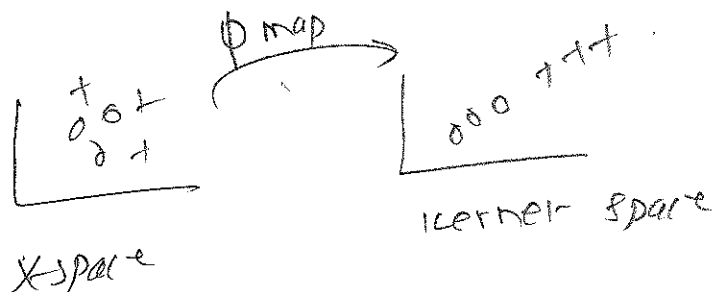
④ cosine similarity

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

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 mutual information

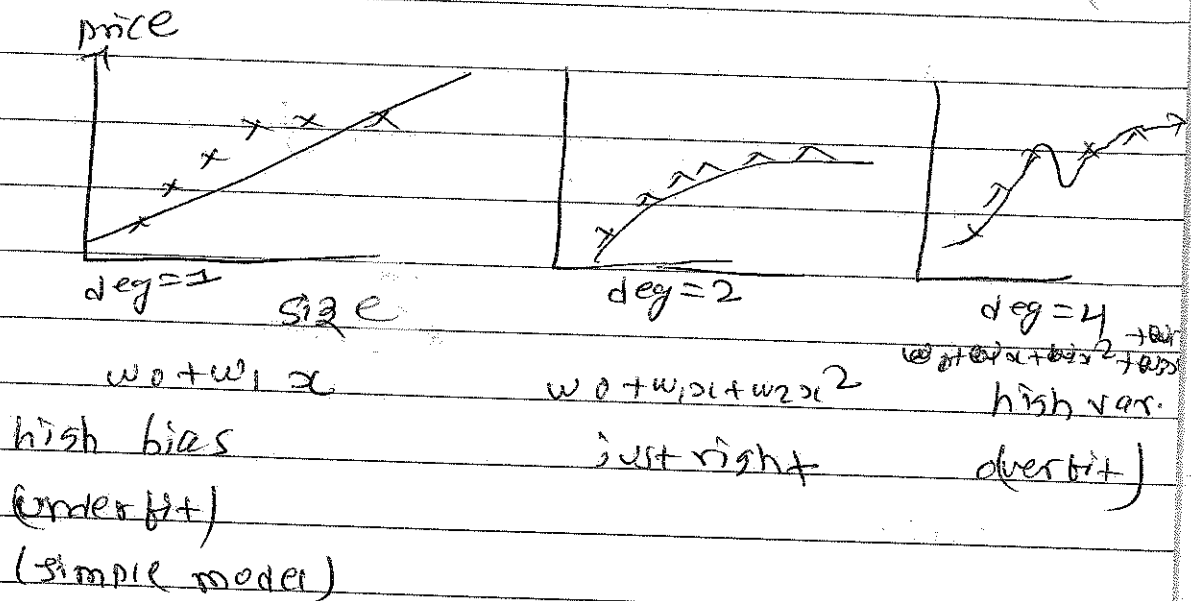
$$\begin{aligned} MI(X, Y) &= \sum_x \sum_y p(x, y) \cdot \ln \left(\frac{p(x, y)}{p(x) p(y)} \right) && = 0 \text{ when } X, Y \text{ indep} \\ &= KL[p(x, y) \parallel p(x) p(y)] && = \max \text{ when } X=Y \end{aligned}$$

bad \rightarrow biased towards high arity features

bad \rightarrow may choose redundant feature

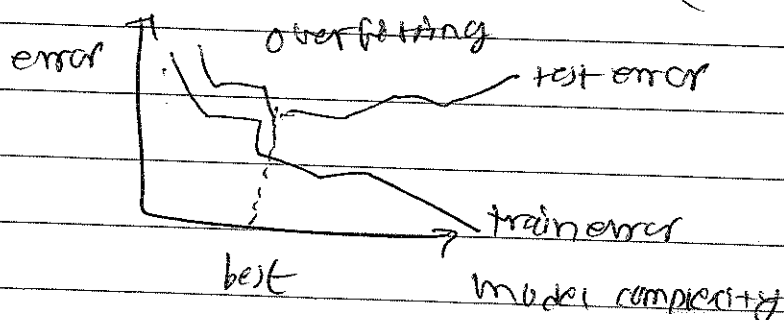
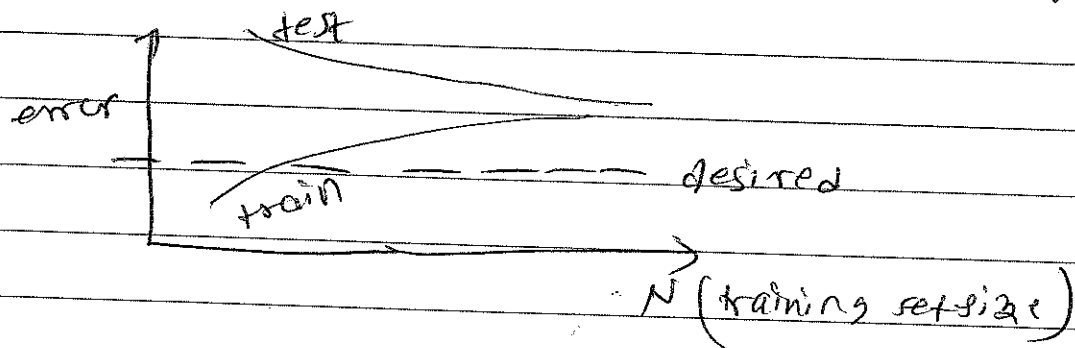
bad \rightarrow works only with nominal features & labels

Bias-variance Trade off

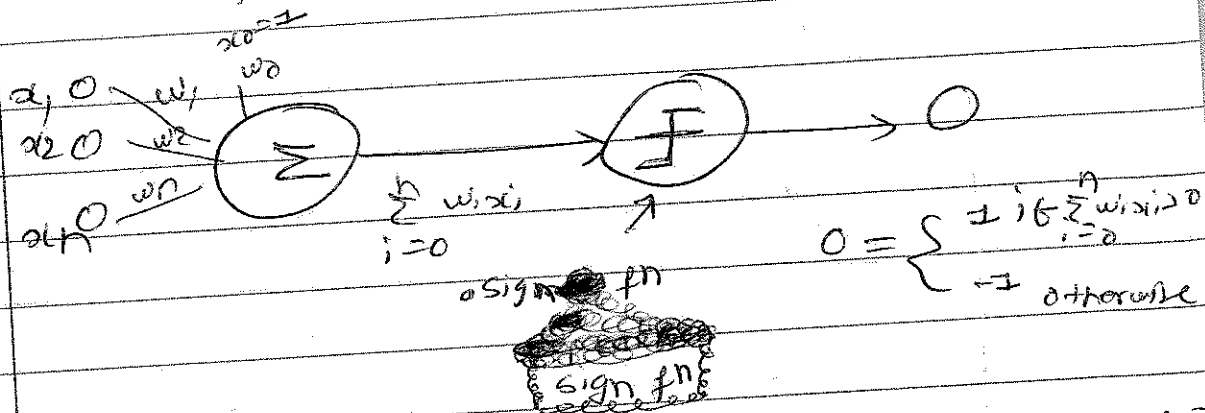


more training eg \rightarrow bias high σ^2 (overfitting)
 smaller feature set \rightarrow bias high σ^2 ("")
 lesser feature set \rightarrow bias high (underfitting)

Learning curve for high bias (underfit)



perceptron



output

$$O = \begin{cases} +1 & \text{if } w_0 + w_1 x_1 + \dots + w_n x_n > 0 \\ -1 & \text{otherwise} \end{cases}$$

$$= \begin{cases} +1 & \text{if } \vec{w} \cdot \vec{x} > 0 \\ -1 & \text{otherwise} \end{cases}$$

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 SVM for regression

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

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

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

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

② SVM for ranking

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

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

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

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

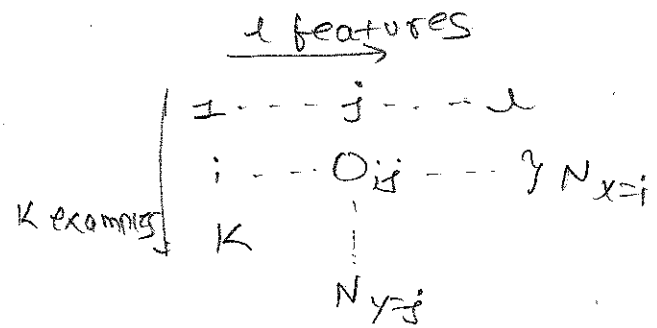
5 mutual information

- works with nominal
- biased towards high parity 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

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

④ SNR

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

binarized

⑤ FTest

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

* 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 $D=2$,

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

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

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

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

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

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

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

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

$$\min E = 0$$

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

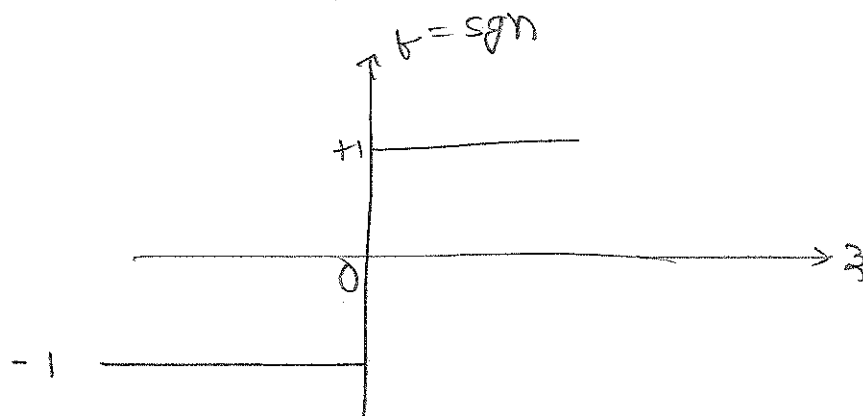
$$0 \xrightarrow{\text{no mistake}} \frac{1}{2} \cdot \frac{N-2}{2} \xrightarrow{\text{N mistakes}} \frac{N}{2}$$

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

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



OP1

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

OP2

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

\Downarrow OP3

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

OP3 has solution \leftarrow

$$\begin{aligned} w_2^* &= w_1^* \gamma \\ b_2^* &= b_1^* \gamma \end{aligned}$$

\Downarrow OP3

rename

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

$$\begin{aligned} w_1 &= w_1^* \gamma \\ b_1 &= b_1^* \gamma \end{aligned}$$

now, decision hyperplanes are,

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

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

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

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

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

③ kmcp kernel multi-class perceptron

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

$$b(x) = w^T x = \sum_n \alpha_n t_n x^T x = \sum_n \alpha_n t_n K(x, x)$$

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

3 for $i = 1, \dots, n$

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

5 if $c_j \neq t_i$ then $(m \neq t_i)$
 $\alpha_{ij} = \alpha_{ij} + 1$ $(d = d + 1)$

Repeat

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

② mcp (concept 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 lap invariant and is the weighted average

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

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

① kcp

define $b(x) = w^T x = \sum_n \alpha_n t_n K(x, x)$

initialize dual parameter $\alpha = 0$

for $i = 1 \dots N$

$h = \text{sign}(b(x_i))$

if $h \neq t_i$ then

$\alpha = \alpha + 1$

Repeat

TEST: $y(x) = \text{sign}(b(x))$

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

expectation & maxing

* Naive Bias

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

• # of parameters = $2^{M+1} = 2^{3+1} = 7$

$p(y=0)$

$P(x_1=1, y=0)$ $P(x_1=1, y=1)$

$P(x_2=1, y=0)$ $P(x_2=1, y=1)$

$P(x_3=1, y=0)$ $P(x_3=1, y=1)$

• # of parameters if no conditional independence = $H(2(2^M - 1)) = H(2(8 - 1))$
 $= H(2(2^3 - 1)) = H(14) = 15$

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 classification maps inputs to discrete outputs
Regression " continuous //

③ PCA \hookleftarrow 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: $w^T x_n \geq 0$ 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

Softmax Regression

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

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

note:

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

$$E_D(w) = \frac{1}{N} \sum_n \sum_k \delta_k(t_n) \ln P(k|x_n)$$

$$= -\frac{1}{N} \sum_n \sum_k \delta_k(t_n) \ln \frac{e^{w_k^T x_n}}{\sum_j e^{w_j^T x_n}}$$

$$= -\frac{1}{N} \sum_n \sum_k \delta_k(t_n) \ln \frac{e^{w_k^T x_n}}{\sum_j e^{w_j^T x_n}}$$

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

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

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

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

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

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

there are K such equations for each classes.

prevent overflow

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

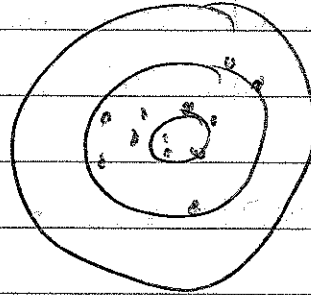
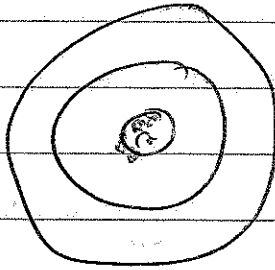
$$c = \max_{1 \leq k \leq K} w_k^T x$$

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

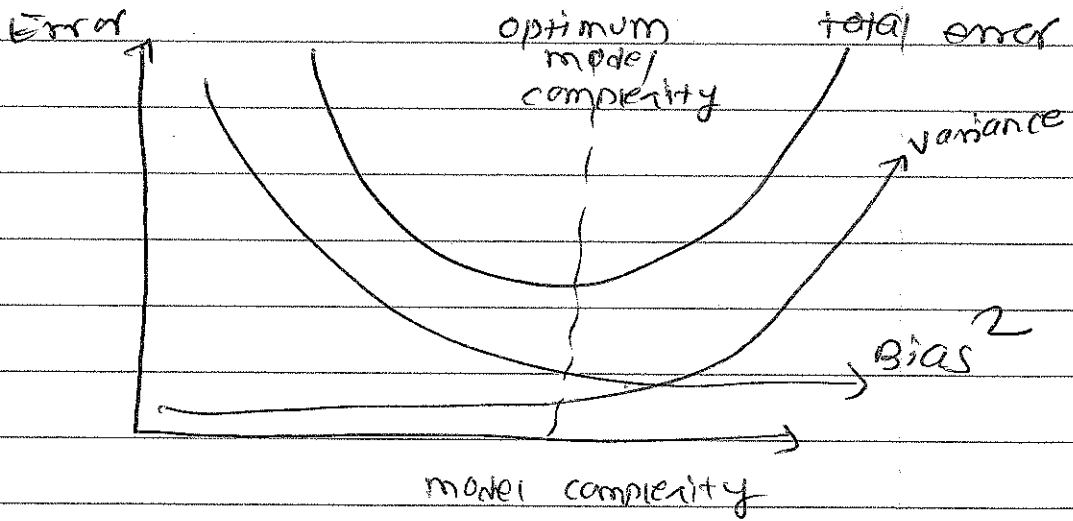
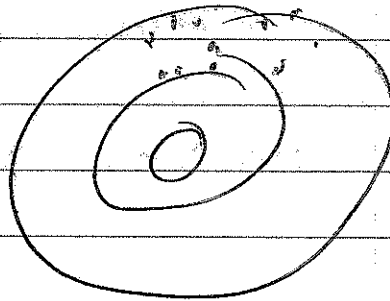
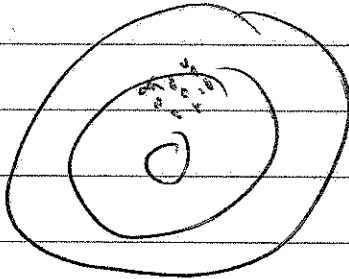
Low variance

High variance

Low Bias



High Bias



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.

MAP solution for LR

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

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

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

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

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

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

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

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

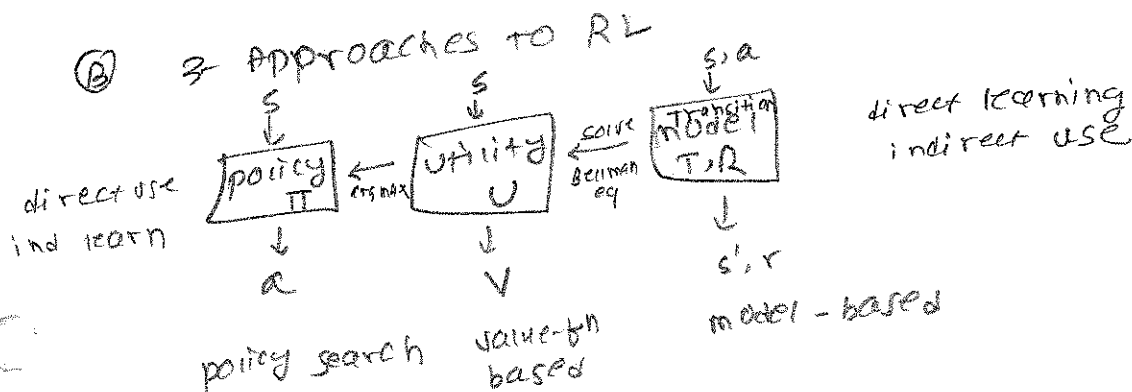
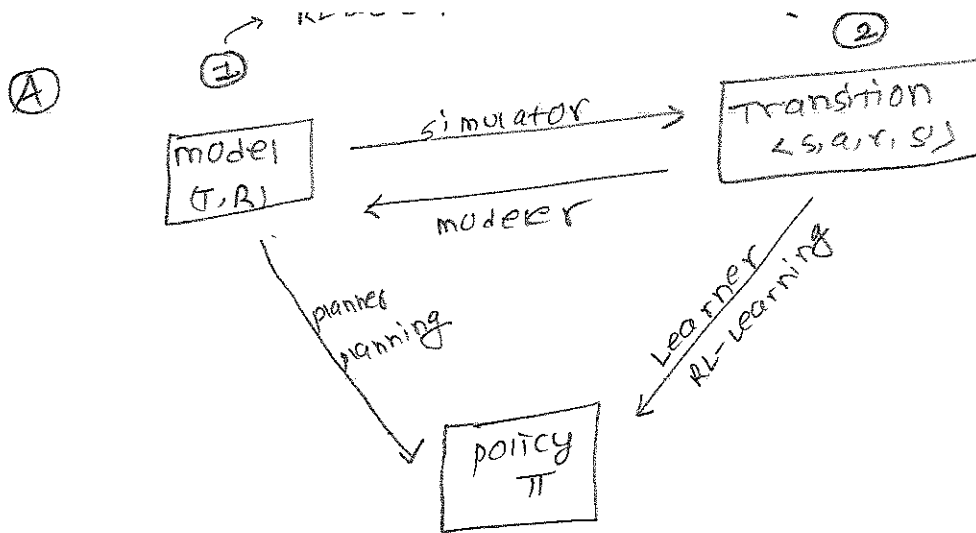
$$= \underset{w}{\operatorname{argmin}} -\sum_n [t_n \ln h_n + (1-t_n) \ln (1-h_n)] + \frac{d}{2} w^T w$$

$$\boxed{\hat{w}_{\text{MAP}} = \underset{w}{\operatorname{argmin}} \underbrace{E_D(w)}_{\substack{\uparrow \\ \text{Ans}}} + \underbrace{E_w(w)}_{\substack{\rightarrow \\ \text{Ans}}}}$$

C

C

C



④ Φ 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 $\Rightarrow U(S) = \max_a Q(S, a)$

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



LINE

EXAMPLE

label

Train	DOC	words (w _i)	class	Total
	3 documents for c1	chinese Beiging chinese chinese chinese Shanghai chinese macao	c1	n1=8
	1 document for c2	Tokyo Japan chinese	c2	n2=3
TEST	5		?	11 words 6 unique
TOTAL		D = D1 + D2 = 3 + 1 = 4		

- chinese 1111
2 Beiging 1
3 Shanghai 1
4 macao 1
5 Tokyo 1
6 Japan 1
- ignore

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

②

category c1 = c

• prior $P(c1) = \frac{|D1|}{|D|} = \frac{3}{4}$

• # of words in class c1, $n1 = 8$

$P(w1|c1) = P(\text{chinese}|c) = \frac{5+1}{8+6} = \frac{6}{14} = \frac{3}{7}$
nk → nk+1 = n1+1 = 5 → Laplace smoothing

$P(w2|c1) = P(\text{Beiging}|c) = \frac{1+1}{8+6} = \frac{2}{14} = \frac{1}{7}$

$P(w3|c1) = P(\text{Shanghai}|c) = \frac{2+1}{8+6} = \frac{3}{14}$

$P(w4|c1) = P(\text{Tokyo}|c) = \frac{0+1}{8+6} = \frac{1}{14}$

$P(w5|c1) = P(\text{Japan}|c) = \frac{0+1}{8+6} = \frac{1}{14}$

category c2 = j (Japan)

• prior $P(c2) = \frac{|D2|}{|D|} = \frac{1}{4}$

• # of words in class c2, $n2 = 3$

$P(w1|c2) = P(\text{chinese}|j) = \frac{1+1}{3+6} = \frac{2}{9}$

$P(w2|c2) = P(\text{Tokyo}|j) = \frac{1+1}{3+6} = \frac{2}{9}$

$P(w3|c2) = P(\text{Japan}|j) = \frac{1+1}{3+6} = \frac{2}{9}$

③ choosing a class

$P(c1|x) \propto P(c1) \prod_i P(w_i|c1)$
prior
 $= \frac{3}{4} \cdot \left(\frac{3}{7}\right)^3 \cdot \frac{1}{14} \cdot \frac{1}{14} \approx 0.0003$

$P(c2|x) \propto P(c2) \prod_i P(w_i|c2)$
prior
 $= \frac{1}{4} \cdot \left(\frac{2}{9}\right)^3 \cdot \frac{2}{9} \cdot \frac{2}{9} \approx 0.0001$

$C^* = \underset{C_k}{\operatorname{argmax}} P(c_k) \prod_{j=1}^n P(w_j|c_k) = \underset{C_k}{\operatorname{argmax}} \{ 0.0003, 0.0001 \}$
 $= \operatorname{Label} 0.0003 = C1$

conditional probabilities

C

C

C