

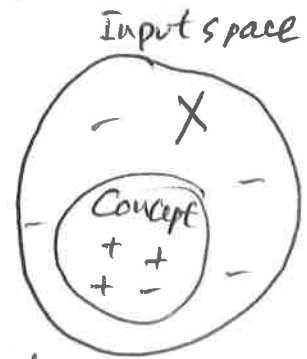
## Lecture - II

## Instance based learning

①

Idea of KNN

- i) Similar objects share same concept
- ii) Classify new examples like similar training examples.



Famous Similarity functions for Real value features:

- 1) Gaussian kernel  $K(x', x) = e^{-(x' - x)^2}$  → depends how far is point  $x$  from  $x'$
- 2) Cosine  $K(x', x) = \cos(x, x')$  → depends on angle of vector

## Lecture - III

## Decision Tree Learning

Example	color	original	presentation	bindee	A+
1	correct (comp, par, guess)				
2					
3					
4					

Instance space = set of all possible objects described by attributes  
 $|X| = 32 \cdot 2 \cdot 2 \cdot 2 = 48$

- \* Target function  $f$ : Maps each instance  $x$  to target label  $y$ . (hidden)
  - \* Hypothesis  $h$ : function that approximates  $f$
  - \* Hypothesis space  $H$ : Set of functions we allow  $H$  approximates  $f$
  - \* Training data  $S$ : Set of instances labeled with target function  $f$ .
- Q: How many functions  $h: X \rightarrow Y$  exist? Ans  $|H| = 2^{48}$

48 instances

(comp, y, y, c, y)	$h_1$	$h_2$	...	$h_n$
	y	y		n
	y	y		n
(guess, no, no, unclear, no)	N	N		n

Just for 48 instances, we have  $2^{48}$  hypothesis functions. Therefore we should restrict hypothesis space.

Example we can restrict "h" that form conjunction p

$$h_p(x) = \begin{cases} \text{yes} & \text{if } p \text{ is true} \\ \text{no} & \text{" " " false} \end{cases}$$

$$P = (\text{correct} = \text{complete}) \wedge (\text{origin} = \text{yes})$$

correct	color	original	presentation	birds
complete	?	yes	?	?

Size of  $|H|$  for conjunctions  $|H| = 4 \cdot 3 \cdot 3 \cdot 3 = 324$   
 $+ 1 = 325$

i.e. we need to ~~search~~ search within 325 hypothesis functions instead of  $2^{48}$ . This is syntactic

#3

## Prediction & overfitting

### Learning Problem

$$P(x, y) = P(y|x)P(x)$$

STEP-1  $P(x)$  is the "world" that produces data instance

Ex 1  $x$  is homework description

$\vec{x}_1 = (\text{complete}, \text{yes}, \text{yes}, \text{clear}, \text{no})$  has some probability

$$P(\vec{x}_1) = 0.2$$

$\vec{x}_2 = (\text{guess}, \text{no}, \text{no}, \text{unclear}, \text{no})$

$$P(\vec{x}_2) = 0.0001$$

Every instance ("48" instances) has some probability. We don't know what is that probability distribution

Ex 2 Patient description: Attributes (fever, age).

Step 2  $P(Y|X)$  is "teacher" that labels  $X$

(2)

(Before we were using function to label " $x_i$ " but now its conditional distribution)

Ex1 Its  $A^+$  homework?

$$P(A^+ | X = \vec{x}_1) = 1 \rightarrow P(\text{not } A^+ | X = \vec{x}_1) = 0$$

Its deterministic labeling.

Ex2

Develops complication in second example?

$$P(\text{complication} | X = (103, 5)) = 0.01$$

$$P(\text{complication} | X = (103, 90)) = 0.3$$

Its not deterministic.

$$P(S) = P(\underbrace{x_1, y_1, x_2, y_2, \dots, x_n, y_n}_{2n \text{ random variables}})$$

IID means :

$$P(S) = \prod_{i=1}^n P(X_i = x_i, Y_i = y_i)$$

i.e. examples are independently drawn i.e. sampling one sample would not effect sampling of other sample.

$$= \prod_{i=1}^n P(X = x_i, Y = y_i)$$

"identically distributed" i.e. we can use same random variables for all samples. because they share same probability distribution.

If something changes overtime, then we will not have IID dataset.

$\text{Err}_P(h)$  is the probability of making an error  $P(h(\vec{x}) \neq y)$  on randomly drawn example from  $P(X, Y)$  (for 0-1 loss)

Overfitting: Hypothesis  $h \in H$  overfits training data  $S$  if there exist  $h' \in H$  such that  $\text{Err}_S(h) < \text{Err}_S(h')$  and  $\text{Err}_P(h) \gg \text{Err}_P(h')$

Occam's Razor: Prefer the simplest hypothesis that fits data.

L#4

## Model Selection And Assessment

(To prevent overfitting, we do model selection)

→ ~~we~~ ~~do~~ learning as prediction

—  $P(x, y)$  is learning task

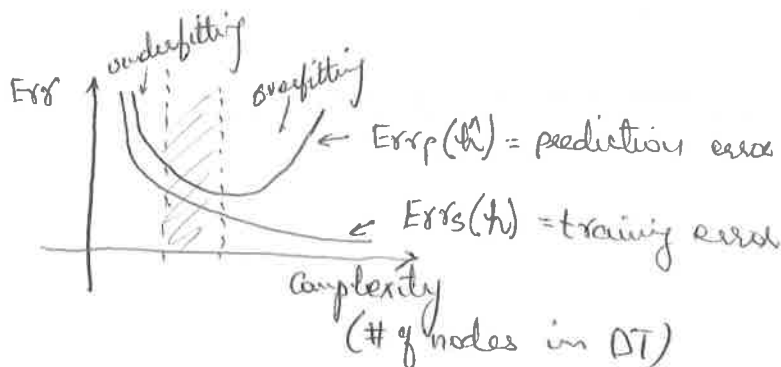
— i.i.d data  $S = (x_1, y_1) \dots (x_n, y_n)$  drawn from  $P(x, y)$

— Sample error  $Err_S(h) = \frac{1}{n} \sum_{(x, y) \in S} \Delta(h(x), y_i)$

— Generalization Error  $Err_P(h) = \sum_{(x, y) \in X} \Delta(h(x), y) P(X=x_i, Y=y_i)$

It is not observed because  $P(X=x_i, Y=y_i)$  is unknown

— Overfitting



In this lecture, we will talk to identify sweet spot in above curve, i.e. in the middle of underfitting and overfitting.

L#5

## Model Selection

Determine the parameters of learning algorithm

— In K.MN, we will find value of  $K$ , similarity measure using model selection

— Decision tree: depth, splitting criterion and many other.

→ In leave-one-out validation, we take  $(n-1)$  examples & test it on the  $n^{th}$  example. we repeat this process " $n$ " times and get unbiased estimator of prediction error. (12:00)

\* What is the prediction error of  $h$ ?

Scenario Is spam filtering rule  $h$  accurate enough ( $Err_p(h) \leq 0.1$ )

→ Binomial D.S

$p = Err_p(h)$   
 $n = \text{no. of ex}$   
 $X = \text{number of errors}$

observation test error  $Err_p(h) = \frac{10}{500} = \frac{x}{n}$

Is there significant evidence that  $Err_p(h) \leq 0.1$ ?

Hypothesis test:

Null hyp  $p \geq 0.1$  (Generalization error  $\geq 10\%$ )

what is the probability to see at most 10 errors under null hypothesis?

OR  $P(X \leq 10 | p = 0.1, n = 500)$

$P(X \leq 10 | 0.1, 500) = 1.1 \times 10^{-12} \leq 0.05$  [This is famous threshold used in p-test]

If null hypothesis is true, then it is extremely unlikely to see ~~few~~ errors less than or equal to 10. Thus we reject null hypothesis.

L#6

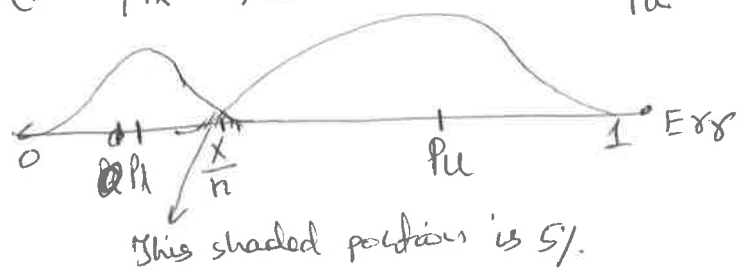
## Linear Classifiers and Perceptrons

Perceptrons = invented by Rosenblatt from Psychology (at the same time in Soviet during 1960's).

What range does the true error rate lie in  $\underbrace{[Err_p(h)]}_{\text{true error rate}}?$   
 → Binomial confidence intervals

Given  $x$  errors on test set, what values  $Err_p(h) \in [p_L, p_U]$  could have plausibly (say 95% confidence) generated  $x$

i.e  $\min_{p_L} P(X \geq x | p_L, n) \leq 0.025$  and  $\max_{p_U} P(X \leq x | p_U, n) \leq 0.025$



$x = \# \text{ of errors}$ , we divide it by  $n$  so that our error ranges from 0 to 1

Example:

→  $n=40$  test examples,  $\text{Err}_{\text{test}}(\hat{h}) = \frac{9}{40}$

→ Normal confidence interval

$$E(X) \in 9 \pm 1.96 \sqrt{\text{Var}(X)}$$

We will use estimated variance

$$\text{Var}(X) = 40 \times \frac{9}{40} \left(1 - \frac{9}{40}\right)$$

$$\text{Var}(X) = 14.176$$

→ 95% normal confidence interval with estimated variance  $\text{Var}(X)$

$$3.824 \leq np \leq 14.176$$

Divided by " $n$ "

$$0.0956 \leq \text{Err}_p(\hat{h}) \leq 0.354$$

After seeing 9 test errors over 40 instances, we can conclude our test error lies in the normal confidence interval  $[0.0956, 0.354]$ .

→ 95% exact binomial confidence interval

$$0.1084 \leq \text{Err}_p(\hat{h}) \leq 0.3844$$

## McNemar Test

	$(x_1, y_1)$	$(x_2, y_2)$	$(x_u, y_u)$	$(x_N, y_N)$
$h_1$	✓	X	✓	X
$h_2$	✓	X	f	✓

$d_1$ : no. of examples  $h_1$  makes error  $h_2$  correct

$d_2$ : " "  $h_2$  " "  $h_1$  "

$\text{Err}_p(\hat{h}_1) = \text{Err}_p(\hat{h}_2) \Rightarrow E(d_1) = E(d_2) \Rightarrow d_1, d_2 \sim \text{Binomial}(p=0.5, n=d_1+d_2)$

Example

1000 test examples

$h_1$  makes 101 errors

$h_2$  makes 111 errors

$$d_1 = 1$$

$$d_2 = 11$$

$$P(D_1 \leq 1 | p=0.5, n=12) = 0.003 \leq 0.5 \text{ (Reject } H_0)$$

(5)

\* Very simple illustration of (Batch) Perceptron Algorithm. L#6 (50:00)  
If the data is linearly separable, we can always find separating hyperplane using perceptron algorithm.

L#7

## Online Learning and Perceptron Mistake Bound

### Worst-case Mistake Bound:

Given set  $C$  of example sequences and learning algorithm " $L$ ", what is the maximum number of mistakes that algorithm " $L$ " makes with hypothesis space " $H$ " makes on any  $S \in C$  before learning a perfect rule?   
 Ans:  $\frac{R^2}{\gamma^2}$   
 classification

→ Implicit Assumption: it exists  $h \in H$  with zero error for all set  $S \in C$ .

### Margin ( $\gamma$ )

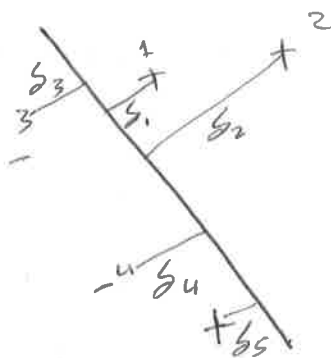
- If margin of example is positive then correct

### Geometric Intuition

Geometric margin is the signed distance of example to hyperplane.

$$d(h_{w,b}, \vec{x}) = \frac{1}{\|w\|} |w \cdot x + b|$$

$$\gamma_{\text{geometric}} = \gamma_{\text{functional}} \cdot \frac{1}{\|w\|}$$



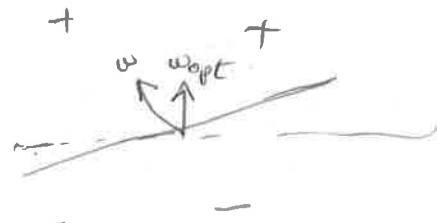
In perceptron, scaling  $\|w\|$  will have no effect while classification,

Unbiased means there is no " $b$ " in  $\|w \cdot x + b\|$ .

To show: perceptron finds hyperplane close to  $w_{opt}$

$\Leftrightarrow$  angle  $(w_{opt}, w_{K+1})$  is small

$\Leftrightarrow \cos(w_{opt}, w_{K+1}) = \frac{w_{opt} \cdot w_{K+1}}{\|w_{K+1}\| + \|w_{opt}\|} \stackrel{!}{=} \text{almost } 1$



$\rightarrow$  find lower bound on  $\cos$

a) lower bound on numerator:  $w_{opt} \cdot w_{K+1}$

b) upper bound on denominator:  $\|w_{K+1}\|$ . we already know that  $\|w_{opt}\| = 1$ .

## Proof by contradiction

If for  $K$ , we can show that lower bound on the  $\cos(w_{opt}, w_{K+1}) > 1$ , then perceptron can not make  $K+1$  errors.

a) Induction:

$$w_{K+1} \cdot w_{opt} \geq f(w_K \cdot w_{opt})$$

$$\begin{aligned} w_{K+1} \cdot w_{opt} &= (w_K + y_i x_i) \cdot w_{opt} \\ &= w_K \cdot w_{opt} + w_{opt} y_i x_i \\ &= w_K \cdot w_{opt} + \underbrace{y_i \cdot x_i \cdot w_{opt}}_{\geq \delta} \\ &\geq w_K \cdot w_{opt} + \delta \end{aligned}$$

~~Overall  $w_{K+1} \cdot w_{opt} \geq w_K \cdot w_{opt} + \delta$~~   
 $\rightarrow w_K \cdot w_{opt} \geq K \cdot \delta$

b) Induction:  $\|w_{K+1}\|^2 \leq g(\|w_K\|^2)$

$$\begin{aligned} \|w_{K+1}\|^2 &= \|w_K + y_i x_i\|^2 \\ &= w_K \cdot w_K + 2 y_i x_i \cdot w_K + y_i x_i \cdot y_i x_i \\ &= w_K \cdot w_K + \underbrace{2 y_i x_i \cdot w_K}_{\leq 0} + \underbrace{x_i \cdot x_i}_{\leq R^2} \\ &\quad \text{b/c we are making update} \\ &\quad \text{e.g. } y_i (w_K \cdot x_i) \leq 0 \\ &\leq w_K \cdot w_K + R^2 \end{aligned}$$



$$\|w_k\|^2 \leq K \cdot R^2$$

(7)

Putting together

$$1 \geq \frac{w_k \cdot w_{opt}}{\|w_k\| \|w_{opt}\|} \geq \frac{S \cdot K}{(\sqrt{R^2 K}) \cdot 1} = \frac{S \sqrt{K}}{R}$$

Solve for K:

$$K \leq \frac{R^2}{S^2}$$

→ Perceptron can not make more than  $\frac{R^2}{S^2}$  updates/errors.

This bound is independent of no. of features.

Remarks

- Perceptron mistake bound is independent of  $N$  (# of features).
- Independent of ordering
- Mistake bound is independent of scaling and rotation of data.
- same for actual behaviour of perceptron
- Downside is that this bound is valid for separable datasets.

L#8

Ensemble methods: Bagging

L#9

Ensemble methods: Boosting

Ensemble method: A class of "meta" learning algorithms because it combines set of classifiers to produce good classification.

\* KNN with  $K=1$  is a learning algorithm because a learning algorithm takes training set as an input and produces  $h(x)$ . If you change training set with  $K=1$ , you will have different decision boundary. Hence ~~KNN~~ KNN with  $K=1$  is learning algorithm.

$$E[X+Y] = E[X] + E[Y]$$

$$E[KX] = K E[X]$$

$$E[XY] = E[X]E[Y] \text{ if } X \perp Y$$

\* Variance:

$$E[(X - E[X])^2] = E[X^2] - (E[X])^2$$

second moment of  $X$

$$\begin{aligned} E[y_i - h_i(x_i)] &= E[y_i^2 - 2y_i h_i(x_i) + h_i^2(x_i)] \\ &= E[y_i^2] - 2E[y_i h_i(x_i)] + E[h_i^2(x_i)] \\ &= y_i^2 - 2y_i E[h_i(x_i)] + E[h_i^2(x_i)] \\ &= \underbrace{(y_i - E[h_i(x_i)])^2}_{\text{variance}} + \underbrace{E[h_i^2(x_i)] - (E[h_i(x_i)])^2}_{\text{variance}} \end{aligned}$$

fixed  
↑  
over subset of training sets

Bagging: You create bag of training set and train set of classifiers over them.

$$h_s = \frac{1}{T} \sum_{i=1}^T \alpha_s h_i(x)$$

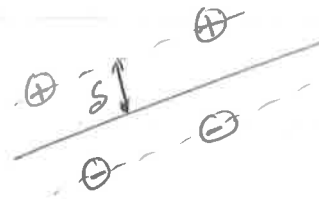
During test time, we ~~average~~ get classification by linear combination  $h_i(x)$ .

L#10 SVM: Optimal hyperplane & soft margin

Compute the optimal hyperplane:

Training sample  $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$

$$\textcircled{1} \max_{(w, b)} \gamma \quad \text{with } \gamma = \min_{(x_i, y_i) \in S} \left\{ \frac{y_i (w \cdot x_i + b)}{\|w\|} \right\}$$



$$\text{Geometric dist} = \frac{y_i}{\|w\|} (w^T x + b)$$

Simplify  $\textcircled{2}$  write min in terms of constraints

$$\max_{w, b} \gamma \quad \text{with } \forall (x_i, y_i) \in S : \frac{y_i}{\|w\|} (w \cdot x_i + b) \geq \gamma$$

$\textcircled{3}$  normalize  $\|w\| = \frac{1}{\gamma}$ , since ~~class~~ classifier is invariant to  $\|w\|$

$$\max_{w, b} \frac{1}{\|w\|} \quad \text{with } \forall (x_i, y_i) \in S \quad \frac{y_i}{\|w\|} (w \cdot x_i + b) \geq \frac{1}{\|w\|}$$

$$\textcircled{4} \quad \text{minimize } \|w\| \quad \text{s.t. } y_i (w \cdot x_i + b) \geq 1$$

$$\textcircled{5} \quad \text{minimize } \|w\|^2 \quad \text{s.t.} \\ \forall (x_i, y_i) \in S \quad y_i (w \cdot x_i + b) \geq 1$$

## 2#10 (continued)

### Soft-Margin SVM

\* Slack variables at solution of QP:

$$\xi_i \geq 1 \Leftrightarrow y_i(w^T x_i + b) \leq 0 \text{ (error)}$$

$$0 \leq \xi_i < 1 \Leftrightarrow 0 < y_i(w^T x_i + b) < 1 \text{ (correctly classified but inside margin)}$$

$$\xi_i = 0 \Leftrightarrow y_i(w^T x_i + b) \geq 1 \text{ (correct with sufficient margin)}$$

Claim  $\sum \xi_i$  is an upper bound to number of training error.

Opt:

$$\frac{1}{2} \|w\|^2 + C \sum \xi_i$$

$$\text{s.t. } y_i(w^T x_i + b) \geq 1 - \xi_i$$

$$\xi_i \geq 0$$

$$\text{HP1: } \text{Error} = 2, \sum \xi_i = 8, \delta = \frac{1}{\|w_1\|} = \frac{1}{\sqrt{2}}$$

$$\text{HP2: } \text{Error} = 0, \sum \xi_i = 0, \delta = \frac{1}{\|w_2\|} = \frac{1}{\sqrt{4}} = \frac{1}{2} = 0.5$$

$$\text{HP3: } \text{Error} = 0, \sum \xi_i = 0, \delta = \frac{1}{\|w_3\|} = \frac{1}{\sqrt{3}}$$

$$\text{HP4: } \text{Error} = 0, \sum \xi_i = 2, \delta = \frac{1}{\|w_4\|} = \frac{1}{\sqrt{2}}$$

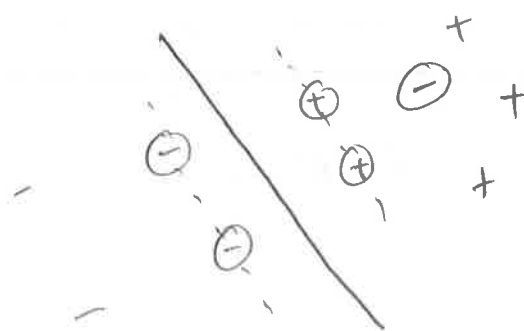
$$\text{HP5: } \text{Error} = 0, \sum \xi_i = 0, \delta = \frac{1}{\|w_5\|} = \frac{1}{\sqrt{2}}$$

$$\text{HP6: } \text{Error} = 0, \sum \xi_i = 0, \delta = \frac{1}{\|w_6\|} = \frac{1}{\sqrt{1.8}}$$

Intuition of SVM Margin It puts weight on those features of  $\vec{w}$  that are decisive for classification.

## #11: Duality and leave-one-out

- \* Support vectors have functional margin  $\leq 1$
- \* Margin is reasonable criterion for assigning weights based on evidence.



All training points in circles are support vectors.

If  $n=10$  in (Batch) Perceptron Algorithm

$$w_{\text{final}} = \alpha_1 y_1 x_1 + \alpha_2 y_2 x_2 + \dots + \alpha_{10} y_{10} x_{10}$$

$\alpha$  = no. of times perceptron makes update over  $x_i, y_i$

$$= \sum_{i=1}^n \alpha_i y_i x_i$$

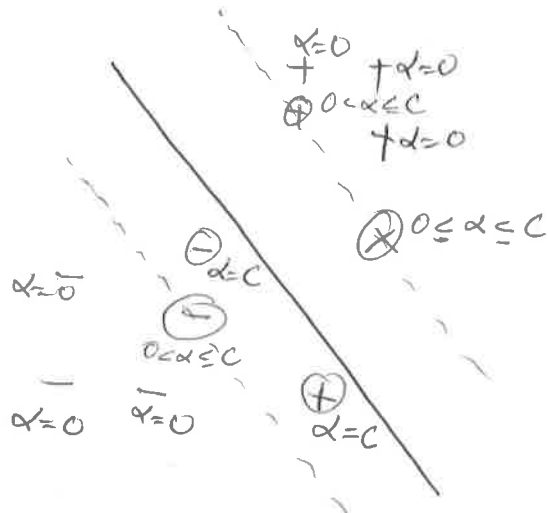
Theorem: The  $w_{\text{final}}$  by perceptron is always the linear combination of training examples

We should check primal & dual perceptron's equivalence

$$\begin{aligned} (w \cdot x_j) &= \left( \sum_{i=1}^n \alpha_i y_i x_i \right) \cdot x_j \\ &= \sum_{i=1}^n \alpha_i y_i (x_i \cdot x_j) \end{aligned}$$

There are two different representation of linear classifier:

- Primal: weight vector " $w$ " & threshold " $b$ "
- Dual: training examples  $(x_1, y_1) \dots (x_n, y_n)$  and dual coefficient  $\alpha_1 \dots \alpha_n$ .



$\alpha=0$  are those training examples that do not disturb decision boundary.

## Dual optimization

— (# of features + 1 parameter + # of examples) free variables in primal

— (# of examples) free variables in the dual.  
     vector of  $\alpha$

— It is convex & QP

— Training examples enter dual only through inner product.

In dual,  $\sum_{i=1}^n \alpha_i y_i = 0$  when we have "b" in primal

Corollary: If  $P(w^*, b^*, \xi^*)$  is solution of primal &  $D(\alpha^*)$  is the solution of dual, then  $D(\alpha^*) = P(w^*, b^*, \xi^*)$

## Leave-one-out error

lets train SVM for whole training set. Then we know:

Lemma

If  $(h(x_i) \neq y_i)$  means  $2\alpha_i R^2 + \xi_i \geq 1$

$\alpha_i, \xi_i$  { from full training set training of SVM  
 $R \geq \|x_i\|$

As a result, we don't need to restrain those examples whos  $2\alpha_i R^2 + \xi_i < 1$ . An immediate case would be for those examples that have  $\alpha=0$ . If we remove those examples, they

decision boundary would be unchanged & we don't need to retrain our classifier.

Theorem: 
$$E_{\text{err}}^{\text{leave}(S)}(\psi_{\text{svm}}) \leq \frac{\# \text{SV}}{\# \text{ of training examples}}$$

i.e. leave one out error is bounded by number of support vectors.

Bound on the expected generalization error

Lemma (unbiased): If  $(x_i, y_i)$  is leave-one-out error then  $\alpha_i R^2 \xi_i \geq 1$

Assumption: 1) we have hard margins  $\rightarrow \xi_i = 0$

If  $\Delta(\psi_i(x_i), y_i) = 1$ , then  $\alpha_i R^2 \geq 1$

$$\Rightarrow E_{\text{err}}^{\text{leave}}(\psi_{\text{svm}}) \leq \frac{1}{n} \underbrace{\#\{\alpha_i R^2 \geq 1\}}_{\text{check it}} \leq \frac{1}{n} \sum_{i=1}^n \alpha_i R^2$$

$\psi_i(x_i)$  = classifier retrained by leaving example "i" out.

$$= \frac{R^2}{n} \sum_{i=1}^n \alpha_i$$

$$\boxed{\sum_{i=1}^n \alpha_i = w^T w} \text{ in Hard margin case}$$

$$= \frac{1}{n} R^2 \cdot \|w\|^2$$

$$= \frac{1}{n} R^2 \cdot \frac{1}{\gamma^2}$$

$\gamma$  = geometric margin

$$= \frac{R^2}{n \gamma^2}$$

(The # of leave-one-out error is upper bounded by  $\frac{R^2}{\gamma^2}$  which is exactly same as perceptron algorithm)

Lemma [Lontz]: The leave-one-out-error is almost an unbiased estimate of the generalization error.

$$E_{\substack{\text{with } (n-1) \\ \text{examples}}} [\text{True error of Alg}] = E_{\substack{\text{size} \\ \text{"n" example}}} \left[ E_{\substack{\text{leave one} \\ \text{out}}} \text{Alg} \right]$$

Theorem If learning task  $P(x, y)$  is separable by unbiased hyperplane, then the expected generalization error of an unbiased hard margin SVM is bounded by

$$E_{\substack{\text{with } (S) = n-1}} [\text{Err}_P(h_{\text{SVM}})] \leq \frac{1}{n} E_{\substack{\text{with} \\ \text{size} = n}} \left[ \frac{R^2}{\bar{\gamma}^2} \right]$$

where  $\bar{\gamma}$  = sample margin of hyperplane computed by SVM.

## L#12 SVM: Kernels

Review:

- Duality for linear classifiers
  - Primal:  $w, b$  ( $N+1$  Params)  $N \rightarrow \# \text{ of features}$
  - Dual:  $\alpha, b$  ( $n+1$  Params)  $n \rightarrow \# \text{ of examples}$
  - $\Leftrightarrow w = \sum \alpha_i y_i x_i$  (connect primal & dual)

- DUAL Perceptron:  $\alpha \in \mathbb{Z}^+$

- DUAL SVM:  $\alpha \in \mathbb{R}$

- Leave-one-out

$$\text{leave-one-out-error} \Rightarrow 2\alpha_i R^2 + \xi_{S_i} \geq 1$$

$$\text{Err}_{100}(h_{\text{SVM}}) \leq \frac{1}{n} \frac{R^2}{\bar{\gamma}^2}$$

$$\text{Err}_{100}(h_{\text{SVM}}) < \frac{\#SV}{n}$$

Mapping  $\phi: X_I \rightarrow X_F$

Input Space:  $X_I$  with dimensionality  $N_I$  (low)

Feature Space:  $X_F$  with  $N_F$  (high)

$\rightarrow a, b, c \text{ (3D)} \rightarrow a, b, c, \underbrace{aa, ab, ac, bb, bc, cc}_{\text{monomials of 2 terms}} \text{ (9D)}$

$W = \begin{pmatrix} -0.7 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix} b = 0$

6-attributes problem on slide 5

Project it  
back  
into  
2-attribute  
space

$h(x) = \text{sign}(-0.7x_1^2 + x_2)$

we will see that we are actually learning non-linear classifier

\* # of monomials of  $N_i$  features and degree less or equal "d" is  $= \binom{N_i + d}{d}$

\*  $\phi(a) \cdot \phi(b) = \phi\left(\begin{pmatrix} a_1 \\ a_2 \end{pmatrix}\right) \cdot \phi\left(\begin{pmatrix} b_1 \\ b_2 \end{pmatrix}\right)$

$= \begin{pmatrix} a_1^2 \\ a_2^2 \\ \sqrt{2}a_1 \\ \sqrt{2}a_2 \\ \sqrt{2}a_1a_2 \\ 1 \end{pmatrix} \cdot \begin{pmatrix} b_1^2 \\ b_2^2 \\ \sqrt{2}b_1 \\ \sqrt{2}b_2 \\ \sqrt{2}b_1b_2 \\ 1 \end{pmatrix}$

$= a_1^2b_1^2 + a_2^2b_2^2 + 2a_1b_1 + 2a_2b_2 + 2a_1a_2b_1b_2 + 1$

$= (a_1b_1)^2 + 2(a_1b_1a_2b_2) + (a_2b_2)^2 + \dots +$

$= (a_1b_1 + a_2b_2 + 1)^2$

$= \left( \begin{pmatrix} a_1 \\ a_2 \end{pmatrix} \cdot \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} + 1 \right)^2$

$= (a \cdot b + 1)^2$

i.e. computing inner product of  $\phi(a) \cdot \phi(b)$  can actually be obtained by  $(a \cdot b + 1)^2$



$$\forall N \nexists d: \phi(a)\phi(b) = (a \cdot b + 1)^d \rightarrow \text{hold for any } d$$

$\uparrow$   
 Kernel:  $K(a, b)$

\* Kernel is like a similarity measure.

### \* SVM with Kernel

How to compute  $b$ ?

- Pick some  $(x_i, y_i)$  with  $0 < \alpha_i < C$

- We know that  $1 = y_i(\omega \cdot \phi(x_i) + b) = y_i \left[ \sum_{j=1}^n \alpha_j y_j K(x_i, x_j) + b \right]$

- Solve for  $b$ :

$$b = y_i - \sum_{j=1}^n \alpha_j y_j K(x_i, x_j)$$

- If we use  $(x_i, y_i)$  whose  $C=1$ , then its margin will be less than 1.

$$\begin{aligned} \|\phi(a) - \phi(b)\|^2 &= (\phi(a) - \phi(b)) \cdot (\phi(a) - \phi(b)) \\ &= \phi(a) \cdot \phi(a) - 2\phi(a) \cdot \phi(b) + \phi(b) \cdot \phi(b) \\ &= K(a, a) - 2K(a, b) + K(b, b) \end{aligned}$$

i.e. Any linear learning algorithm that I can express in terms of distances or in terms of inner products, I can just plug in kernel and make it non-linear.

\* Any boolean function can be represented by ~~bool~~ kernels as well.

L#13

## learning ranking functions with SVMs

Why kernel?

- make linear learner non-linear
- make linear learner with non-vectorial data

What is kernel?

- It's a kernel correspond to inner product in feature space  
i.e.  $K(a, b) = \phi(a) \cdot \phi(b)$
- Kernel is efficient to compute even if  $\phi(a)$  is inefficient.

Constructing kernel?

- Construct  $\phi(a)$ , then find kernel  $K(a, b)$ .
- Check that Gram matrix is always PSD & symmetric.
- Construct it from simple kernels.

"Kernelize" a learning algorithm.

- replace inner product by kernel.

\* For ranking, we learn utility function.

L#14

## Discriminative Vs Generative learning

Course So far

- Decision Trees learning
  - ~~greedy~~ greedily construct small tree with low training error
  - Problem: no theory
- Perceptron
  - find separating hyperplane
  - mistake bound
  - Problem: should be no noise
- SVM
  - soft margin
  - theory about prediction error
  - non linear / non-vectorial data via kernels.
  - learning to rank

Uptill now, we have hypothesis space & we are searching for hypothesis that have smallest training error (empirical risk minimization) or more generally discriminative method.

### Discriminative training of linear rules

$$\min_{w, b} \underbrace{R(w)}_{\text{Regularizer}} + \underbrace{C \sum_{i=1}^n \Delta(w^T x_i + b, y_i)}_{\text{training loss / empirical risk / Training error}}$$

Regularization Parameter

#### • Soft Margin SVM

$$R(w) = \frac{1}{2} \|w\|^2 = 0.5 w^T w$$

$$\Delta(\bar{y}, y_i) = \max(0, 1 - y_i \bar{y}) \text{ [margin=1]}$$

#### • Perceptron

$$R(w) = 0$$

$$\Delta(\bar{y}, y_i) = \max(0, -y_i \bar{y}) \text{ [margin=0]}$$

#### • Linear Regression

$$R(w) = 0$$

$$\Delta(\bar{y}, y_i) = (y_i - \bar{y})^2$$

#### • Ridge Regression

$$R(w) = \frac{1}{2} w^T w$$

$$\Delta(\bar{y}, y_i) = (y_i - \bar{y})^2$$

#### • Lasso

$$R(w) = \frac{1}{2} \sum |w_i|$$

$$\Delta(\bar{y}, y_i) = (y_i - \bar{y})^2$$

#### • Regularized Logistic Regression / Conditional Random Field

$$R(w) = \frac{1}{2} w^T w$$

$$\Delta(\bar{y}, y_i) = \log(1 + e^{-y_i \bar{y}})$$

#### • SVM Regression

$$R(w) = \frac{1}{2} w^T w$$

$$\Delta(\bar{y}, y_i) = \max(0, \xi - y_i \bar{y})$$


Continue from 31:00

Example:  $X$  is patient in hospital

$$P(y=1|X=x) = 0.7$$

$$P(y=-1|X=x) = 0.3$$

→ Predict +1 has error of 0.3

→ If we predict -1, we will have an error of 0.7.

Bayes Rule:  $\hat{y}_{\text{bayes}}(\vec{x}) = \underset{y \in Y}{\operatorname{argmax}} P(Y=y|X=\vec{x})$

## Bayes Error Rate

What is  $\text{Err}_p(h_{\text{Bayes}})$ ?

- Given an instance  $x$ ,  $h_{\text{Bayes}}$  has prob of making an error of  $\min_y [1 - P(Y=y|X=x)]$ .

- Average over all  $x$ :

$$\text{Err}_p(h_{\text{Bayes}}) = \sum_{\substack{X \in \mathcal{X} \\ \text{known}}} P(X=x) \min_y (1 - P(Y=y|X=x))$$

This is Bayes optimal error rate.

As  $P(X,Y) = P(Y|X)P(X)$ : we know that  $P(X,Y)$  is our learning task and we ~~can't~~ don't know <sup>this</sup> distribution. For instance in ~~text~~ text classification task,  $P(X,Y)$  is the distribution of authors writing documents.

## Generative vs Discriminative Models

### Process:

- Generator: Generate descriptions according to distribution  $P(X)$
- Teacher: Assigns a value to each distribution based on  $P(Y|X)$

↓  
Training Examples  $(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n) \sim P(X,Y)$

### Discriminative Model

- Select classification rules  $H$  to consider (Hypothesis space)
- find  $h$  from  $H$  with lowest training error
- Argument: low training error will lead to low prediction error
- Examples: SVM, decision trees, Perceptron

### Generative Model

- select set of distributions to consider for modelling  $P(X,Y)$
- find distribution that matches  $P(X,Y)$  on training data
- Argument: If match close enough, we can use Bayes' decision rule
- Example: naive Bayes, HMM

L#4  
 Deriving a generative learning algorithm

$$h(x) = \underset{y}{\operatorname{argmax}} [P(Y=y|X=x)]$$

Binary Classification

$$P(Y=1|X=x) = \frac{P(X=x|Y=1)P(Y=1)}{P(X=x)}$$

$$P(Y=-1|X=x) = \frac{P(X=x|Y=-1)P(Y=-1)}{P(X=x)}$$

$$\begin{aligned} \Rightarrow h(x) &= \underset{y}{\operatorname{argmax}} (P(Y=y|X=\vec{x})) \\ &= \underset{y}{\operatorname{argmax}} \left[ \frac{P(X=\vec{x}|Y=y)P(y)}{P(X)} \right] \end{aligned}$$

As we want max over  $y$ , we can remove  $P(X)$

$$= \underset{y}{\operatorname{argmax}} P(X=\vec{x}|Y=y)P(y)$$

Now we have our model as follows:

$$h(x) = \underset{y}{\operatorname{argmax}} P(X=\vec{x}|Y=y)P(y)$$

Intuition: Let's say we have two bags of +ve & -ve classes. We flip a coin & see which class is this. It will be  $P(y)$ .

Then we go to respective bag & sample  $\vec{x}$ . It is  $P(X|Y)$ .

In this way, we generate our dataset. We sample our class based on  $P(Y)$  & then we sample  $X$  based on  $P(X|Y)$ . So if we know two distributions i.e.  $P(Y)$  &  $P(X|Y)$  then we can classify in generative way.

Intuition

$P(y)$  = How likely is class " $y$ " occurring a priori

~~$P(X)$~~   $P(X=\vec{x}|y=+1)$  = How likely do I see  $\vec{x}$  in class +1

$P(X=\vec{x}|y=-1)$  = How likely do I see  $\vec{x}$  in class -1

## Naive Bayes Assumption

$$P(X=\vec{x} | Y=1) = P(X_1=x_1, X_2=x_2, \dots, X_n=x_n | Y=1)$$

$$= \cancel{P(X_1=x_1 | Y=1)} P(X_2=x_2 | Y=1) \dots P(X_n=x_n | Y=1)$$

$$P(X=\vec{x} | Y=-1) = \dots$$

i.e. all features are independent

## Example on slides on NB

$$\hat{P}(Y=1) = 3/4 \quad \hat{P}(Y=-1) = 1/4$$

$$\hat{P}(x_{\text{fever}} = \text{high} | Y=1) = 2/3, \quad \hat{P}(x_{\text{fever}} = \text{low} | Y=1) = 1/3, \quad \hat{P}(x_{\text{fever}} = n | Y=1) = 0$$

$$\hat{P}(x_{\text{fever}} = \text{high} | Y=-1) = 0, \quad \hat{P}(x_{\text{fever}} = \text{low} | Y=-1) = 1, \quad \hat{P}(x_{\text{fever}} = n | Y=-1) = 0$$

## Classify new example

$$P(Y=1 | X) = P(Y=1) \prod P(X_i=x_i | Y=1)$$

$$= 3/4 \left( \underset{\substack{\uparrow \\ \text{fever}=\text{hi}}}{3/4} \times \underset{\substack{\downarrow \\ \text{cough}=\text{no}}}{1/3} \times \underset{\substack{\downarrow \\ \text{toes}=\text{yes}}}{2/3} \right)$$

$$= 1/11 \text{ or check it}$$

$$P(Y=-1 | X) = P(Y=-1) \prod P(X_i=x_i | Y=-1)$$

$$= 1/4 \times 0$$

$$= 0 \text{ (use Laplace smoothing)}$$

we now need to estimate polynomial number of parameters  
i.e. # of params = # of features  $\times$  # of classes + # of classes  
Before making naive Bayes assumption we had exponential # of parameters

Now look at:

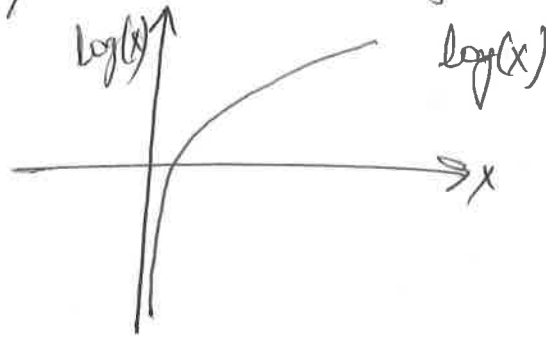
$$h_{\text{naive}}(\vec{x}) = \underset{y \in \{+1, -1\}}{\operatorname{argmax}} \left\{ P(Y=y) \prod_{i=1}^n P(X_i=x_i | Y=y) \right\}$$

Apply log & it will not change max.

$$h(x) = \underset{y}{\operatorname{argmax}} \left[ \underbrace{-\log P(Y=y)}_{\text{by}} + \underbrace{\sum_{i=1}^n \log P(X_i=x_i | Y=y)}_{W_y \cdot x} \right]$$

i.e. it will become a linear ~~rule~~ function. or in naive Bayes we actually learn linear rule in log scale.

we take negative log ~~base~~ (see graph). b/c log of real number b/w 0 & 1 is negative.



L#15

## Generative Model for Classification

Review:

- Discriminative learning (For vs Expected Risk Minimization)
  - Find rule with low training error
- Generative learning
  - Estimate  $P(X, Y)$  from our training data and then derive classification rule via Bayes decision rule.
- Bayes Decision rule:  $h(x) = \underset{y}{\operatorname{argmax}} \{P(Y=y|X=x)\}$ 
  - If  $P(X, Y)$  is known, then Bayes decision rule is optimal
- Multivariate Naive Bayes algorithm
 

$$h(x) = \underset{y}{\operatorname{argmax}} \left\{ \frac{P(X|Y) P(Y)}{P(X)} \right\}$$

$$= \underset{y}{\operatorname{argmax}} \{ P(X|Y) P(Y) \}$$

$P(Y)$  = prior prob of class before looking at  $X$ .

Assume:

$$P(X=x|Y=y) = \prod_{i=1}^N P(X_i=x_i|Y=y)$$

In Generative Model, we need to estimate following two distributions from ~~lower~~ training data: (i)  $P(Y=y)$  (ii)  $P(X_i=x_i|Y=y)$

\* For different generative algorithms, we have different assumptions over class conditional model i.e.  $P(X|Y)$ . In naive Bayes, each feature of  $X$  is independent  $\&$  in HMM,  $x_t$  depends on  $x_{t-1}$ .

\* Why do we call it generative ~~prob~~ approach?

Because we are modelling generative process that generate our training data.

Consider classification rule (Bayes classification rule):

$$h(x) = \underset{y}{\operatorname{argmax}} \underbrace{P(X=\vec{x} | Y=y)}_{\text{class conditional}} \underbrace{P(Y)}_{\text{prior}}$$

Here we assume that our ~~data~~ can be decomposed as follows

Wagdy K

$$P(X=\vec{x} | Y=y) = \prod_{i=1}^N P(X_i=x_i | Y=y)$$

### ASSUMPTION of DATA Generation

Now imagine we have " $N$ " bags for +ve &  $N$  bags for -ve training examples

for -ve training examples  
BINARY CLASSIFICATION

BINARY CLASSIFICATION

Imagine we have "N" bags for "N" features that belong to +ve class

" " " " " " " " " " -ve "

We flip a coin, and it turns +ve. We will go to 1st bag that belong to +ve example & ~~stop~~ pull out feature value. Then will move to 2nd, ...,  $N^{th}$  bag & we will form  $\vec{x_1} = (x_1, x_2, \dots, x_N)$  1st example labelled as +ve.

In the same way, we will generate whole training set.  
So, we assumed that our data was generated in the above explained way & we therefore learnt two distributions.

explained way & we therefore must  
 $P(X=x | Y=y)$  &  $P(Y=y)$  from data. That's why ~~we~~ we  
 call this approach, generative approach. Because we  
 just modelled generative process.

just modelled generative process.  
The above model is convenient b/c we had discrete variables i.e.  $\vec{X}_i = (x_1, x_2, \dots, x_N)$  takes discrete ~~variables~~ values.



Now consider different problem:

How would we construct for generative model for this data?

Ans: Linear Discriminant Analysis

$$h(x) = \underset{y}{\operatorname{argmax}} \{ P(X=x | Y=y) P(Y=y) \}$$



Now we need to model:

- i)  $P(Y=y)$  that's easy. Just count +ve & -ve classes  
 ii)  $P(X=x | Y=y)$ : what should be the assumption for this distribution? Ans: Gaussian or spherical distribution. Fit a gaussian.

Our Assumption:

$$P(X=\vec{x} | Y=+1) = N(\vec{\mu}_+, \Sigma_+) = \frac{1}{2\pi^{1/2} |\Sigma_+|^{1/2}} e^{-\frac{1}{2}(\vec{x}-\vec{\mu}_+)^T \Sigma_+^{-1} (\vec{x}-\vec{\mu}_+)}$$

$$P(X=\vec{x} | Y=-1) = N(\vec{\mu}_-, \Sigma_-)$$



Since we have round shapes

$$\Sigma_+ = \Sigma_- = \Sigma = \frac{1}{2} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \text{ (Spherical Gaussian with equal covariance)}$$

If our data is elongated, we can normalize data with  $\mu=0$  &  $\Sigma=1$  which will give above illustration of data.

$$h(x) = \underset{y}{\operatorname{argmax}} \{ N(\mu_y, \Sigma_y) P(Y=y) \}$$

Dropping out constants

$$= \underset{y}{\operatorname{argmax}} \{ e^{-\frac{1}{2}(\vec{x}-\mu_y)^2} P(Y=y) \}$$

Apply log

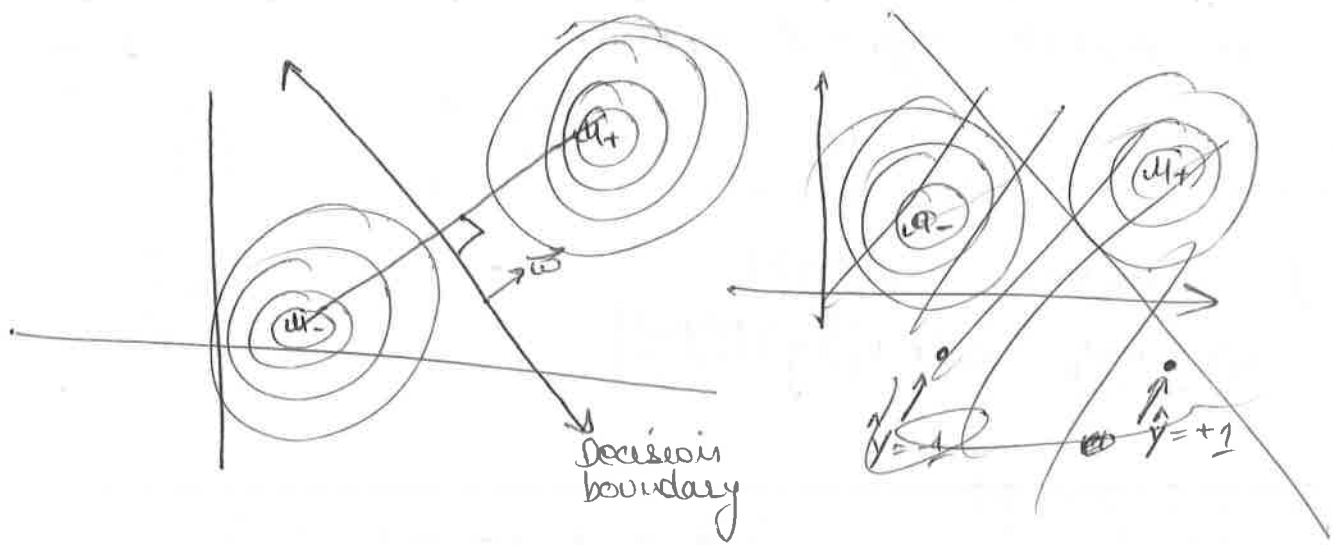
$$= \underset{y}{\operatorname{argmax}} \left\{ -\frac{1}{2}(\vec{x}-\mu_y)^2 \log P(Y=y) \right\}$$

Assume  $P(Y=+1) = P(Y=-1) = 0.5$

$$= \underset{y}{\operatorname{argmin}} \left\{ -\frac{1}{2}(\vec{x}-\mu_y)^2 \right\}$$

$$= \underset{y}{\operatorname{argmin}} \{ \|\vec{x} - \mu_y\| \}$$

i.e we can classify by measuring the distance b/w  $\vec{x}$  & mean



If we assume that prior probabilities are same then we actually classify based on distances. If we ~~we~~ classify based on distances, then we again learnt linear decision boundary as shown in above fig. That decision boundary is  $\perp$  to the line passing through  $u_+$  &  $u_-$ . But we have to assume that we have unit covariance for both classes.

$$w = u_+ - u_-$$

$$b = ? ?$$

We got linear classifier with minimization error framework.

## Naive Bayes Classifier (Multinomial)

~~Assume~~

Bayes Rule:

$$h(x) = \underset{y}{\operatorname{argmax}} P(X=x | Y=y) P(Y=y)$$

$$\underbrace{P(X=x | Y=y)}_{\text{How to model it}} P(Y=y)$$

b/c  $X$  = sequence of words

Assume:

Documents are generated by independently drawing " $l$ " words according to a class conditional distribution.

$$P(X=x | Y=1) = P(w_1=w_1, w_2=w_2, \dots, w_l=w_l | y=1)$$

$$= \prod_{i=1}^l P(w_i=w_i | y=1)$$

$$P(\vec{x} = \vec{x} | y = -1) = \dots$$

Now, let's see its generative process & consider the following example.

$$x_1 = (\text{The, art, of, Programming}) \quad y_1 = +1$$

$$l_1 = 4$$

Now we have two bags of ~~comp~~ words. One bag contains words for computer science documents and another bag contains words for the articles that do not belong to computer science. We flip a coin ~~at each step~~ & get class label. Then we go to the respective bag & draw each word independently of others i.e. by replacement. So for  $l_1 = 4$ , we will draw 4 words & during each draw we put the previously drawn word back to the bag. This is different from ~~multinomial~~ multivariate Naive Bayes model where we had "N" different bags for "N" features and for each class.

The key observation here is that we are not modelling syntax. We are just drawing words whose order does not matter. We could have used multivariate ~~Naive Bayes~~ <sup>method</sup> ~~assumption~~ here for modeling  $P(\vec{x} | y)$  but then we would have to go through huge amount of words. Therefore, we just modeled those words which are present in example. How to ~~model~~ estimate  $P(w = w | y = y)$ , see slides. It's also linear decision rule. (See homework)

- \* We know text is not generated in the above described way. Therefore this model is gross-over-simplification. But this model is superfast.
- \* Statistically, order & order does not matter.
- \* LDA can be kernelized.
- \* Good slides for comparing different model at the end of this lecture.

L#16

# Modeling Sequence Data: Markov Models

## Review

- Generative Modeling:  $h(x) = \underset{y \in Y}{\operatorname{argmax}} \{ \underbrace{P(x|y)}_{\text{Class conditional distribution}}, \underbrace{P(y)}_{\text{Prior: How likely is "y" before seeing data}} \}$

- Multinomial Naive Bayes (applied to sequence, documents)

Assumption:  $P(X=x|Y=y) = \prod_{i=1}^l P(w_i=w_i|Y=y)$  (bag of words)

Estimate:  $\underbrace{P(w|y)}_{\substack{\# \text{ of words} \\ \times \# \text{ of classes}}}, \underbrace{P(y)}_{\# \text{ of classes}}$

- Linear Discriminant Analysis:

Assume:  $P(X=\vec{x}|Y=y) = N(\vec{\mu}_x, \Sigma)$

Estimate:  $\vec{\mu}_y$

## Less Naive Bayes

Assume:

$P(X=x|Y=y) = P(w_1, w_2, \dots, w_l|y)$  {Here  $X$  is a sequence}

(Rule of Probability)  $= P(w_1|y) P(w_2=w_2|w_1=w_1, y=y) \dots P(w_l=w_l|w_1=w_1, \dots, w_{l-1}=w_{l-1}, y=y)$   
 $= P(w_1=w_1|Y=y) P(w_2=w_2|w_1=w_1, Y=y) P(w_3=w_3|w_1=w_1, w_2=w_2, Y=y) \dots P(w_l=w_l|w_1=w_1, \dots, w_{l-1}=w_{l-1}, Y=y)$

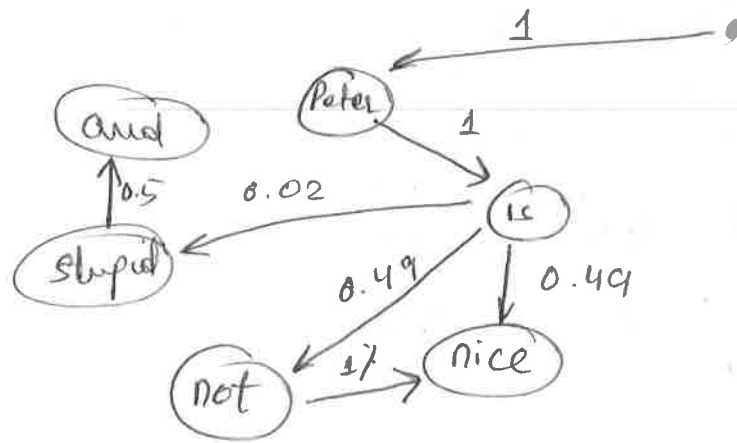
## Generative Model



← (Model for insulting sentence)

\* Edges represent class-conditional probability

## NOT INSULT



(sum of edge from node = 1)

We flip a coin and go to respective graph. Then we start from ~~not~~ start node and perform random walk "1" times on respective graph to generate sentences of length "1".

Assume:  $P(y) = 0.5$  (uniform)

$$\begin{aligned} P(y=1(\text{insulting})) P(X=x_1 | y=\text{insult}) \\ = 0.5 P(\text{Peter, is, nice, and, not, stupid} | \text{insult}) \\ = 0.5 \times 1 \times 1 \times 0.02 \times 0.5 \times 0.5 \times 0.01 \end{aligned}$$

$$P(y=-1) P(X=x_2 | y=-1) = 0.5 \times \dots$$

$$\# \text{ of params} = [(\# \text{ of words})(\# \text{ of words}) - 1] \times \text{classes}$$

Generative Model for POS tagging

$$h(x) = \arg\max_{y \in Y} \{ P(X=x | Y=y) P(Y=y) \}$$

This was "2" classes when there was no sequence. With sequence its exponential

Prior:  $P(Y=(y_1, \dots, y_n))$

- Dependencies follow Markov Model

$$P(Y=(y_1, \dots, y_n)) = \underbrace{P(Y_1=y_1)}_{\text{start prob}} \prod_{i=2}^n \underbrace{P(Y_i=y_i | Y_{i-1}=y_{i-1})}_{\text{transition}}$$

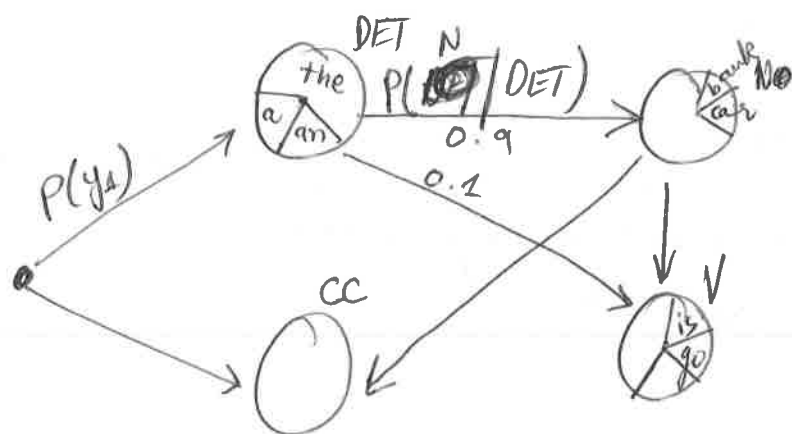
## Class - Conditional Model $P(X=x|Y=y)$

- Assume that word  $x_i$  is independent of
  - all other words
  - all tags except  $y_i$

$$P(X=(x_1, \dots, x_L) | Y=(y_1, \dots, y_L)) \\ = \prod_{i=1}^L P(X_i=x_i | Y_i=y_i)$$

i.e. we have markov model over tags & simple model over  $P(X|Y)$ .

## Generative Process



\* Circles are bag of words with state (DET, N, CC, V)

SAMPLE  $Y$  & then  $X$ :

$y = (\text{DET}, \text{random walk}) \rightarrow y = (\text{DET}, \text{noon}) \rightarrow y = (\text{DET}, \text{N}, \text{V})$   
 $x = (\text{The}, \text{draw brown bag of words}) \rightarrow x = (\text{She}, \text{car}) \rightarrow x = (\text{The}, \text{car}, \text{is})$

→ we will have a random walk over above graph which will be guided by  $P(y_i|y_{i-1})$ . We will generate label sequence & then we will pick a word from respective <sup>bag</sup> ~~word~~. We will sample words based on their frequencies/probabilities inside bag.

\* We need fast algorithms to compute argmax ~~for~~ for sequence of labels.

L#17

## Modeling Sequence Data: HMMs and Viterbi

### Review:

#### Prob. Dist over sequences

##### - Markov Model

##### - Prob distribution over sequences

$$P(S = (s_1, \dots, s_l)) = \underbrace{P(S_1 = s_1)}_{\text{Start Prob}} \prod_{i=2}^l \underbrace{P(S_i = s_i | S_{i-1} = s_{i-1})}_{\text{transition Prob}}$$

- Sampling from above distribution is random walk over graph whose edges are transition prob

- # of sequences of length " $l$ " with " $K$ " symbols is  $K^l$ . We have defined probability distribution over  $K^l$  sequences by Markov assumption. Its parameters are  $K^2 + K$

$\uparrow$                        $\uparrow$   
# of trans              # of start prob

##### - Less Naive Bayes Classifier

- Assumption: Class conditional model  $P(X|Y)$  are Markov Models

- Till here, we were predicting binary variables

##### - Sequence Prediction:

- Application: POS tagging

- TAGS WAS: Given sequence  $X$ , predict sequence  $Y$

- Generative Model:

$$P(X=x|Y=y) \times P(Y=y)$$

— In  $P(X=x|Y=y) P(Y=y)$

Since it is huge as  $\vec{x}$  is a sequence of words, it will have huge # of conditional distributions. So, we made an assumption as follows:

$$P(X=x|Y=y) = \prod_{i=1}^L P(X_i=x_i|Y_i=y_i) \text{ i.e. we draw } x_1$$

$$\vec{x} = ( \underset{x_1}{I}, \underset{x_2}{am}, \underset{x_3}{living}, \underset{x_4}{in}, \underset{x_5}{PAK} )$$

by sampling  $y_1$  & then we go to bag of words ~~label~~ belonging to state  $y_1$ . Then we sampled  $x_1$ .

$$P(Y=y) = P(Y_1=y_1) \prod_{i=2}^L P(Y_i=y_i|Y_{i-1}=y_{i-1}) \rightarrow \text{Modeling POS sequence by random walk.}$$

— Put everything together, we will see HMM.

$$\begin{aligned} \text{i.e. } P(X,Y) &= P(X|Y) \times P(Y) \\ &= P(x_1, \dots, x_L, y_1, \dots, y_L) \end{aligned}$$

— See Viterbi Example in slide #3.

Example:

$$\begin{aligned} P(X=(I, bank, at, CFCU), Y=(N, V, N, V)) &= ?? \text{ what is the probability of this sequence} \\ &= P(Y_1=N) P(X_1=I|Y_1=N) \cdot P(Y_2=V|Y_1=N) P(X_2=bank|Y_2=V) \\ &\quad P(Y_3=N|Y_2=V) P(X_3=at|Y_3=N) \cdot P(Y_4=V|Y_3=N) P(X_4=CFCU|Y_4=V) \\ &= \text{Get value from slides} \\ &= (0.1 \times 0.01) \cdot (0.19 \times 0.4) \times (0.3 \times 0.01) \times (0.19 \times 0.01) \\ &= 4.332 \times 10^{-10} \text{ (very unlikely)} \end{aligned}$$

— For HMM, we need if we have "K" symbols then we need to estimate following params:

$$\begin{aligned} \text{trans} &= K^2 \\ \text{emission} &= L \text{ (# of distinct words in English language or vocabulary)} \end{aligned}$$



\* If we don't simplify our model, then number of parameters would be exponentially huge.

\* Viterbi is general case of Belief Propagation.

### \* Viterbi Algorithm

\* Prob of most likely sequence of states/tags ending in <sup>state "y" at</sup> position "i".

$$S_y(i) = \max_{(y_1 \dots y_{i-1}) \in S^{i-1}} P(X_1=x_1 \dots X_i=x_i, Y_1=y_1, \dots, Y_{i-1}=y_{i-1}, Y_i=y)$$

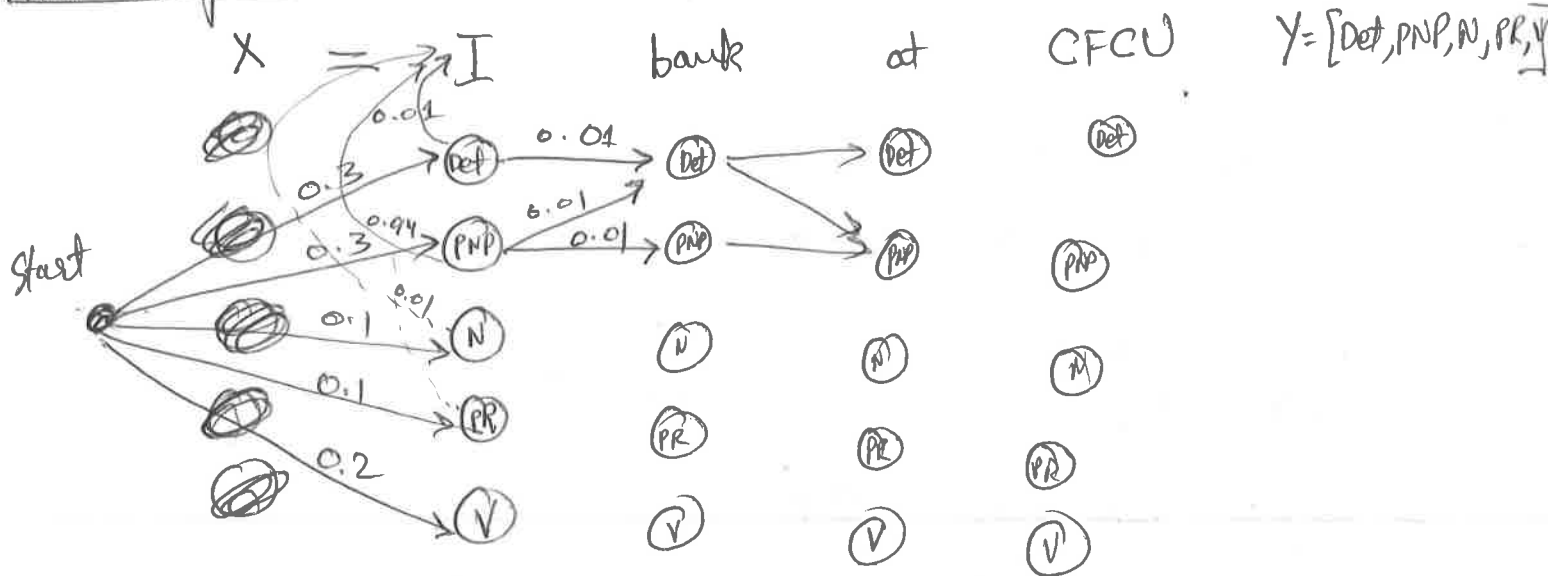
$$= \max_{(y_1 \dots y_{i-1}) \in S^{i-1}} \left\{ P(X_1=x_1 | Y_1=y_1) P(y_1) \prod_{j=2}^{i-1} P(X_j=x_j | Y_j=y_j) P(y_j | Y_{j-1}=y_{j-1}) \right.$$

$$\left. P(X_i=x_i | Y_i=y) P(y_i=y | Y_{i-1}=y_{i-1}) \right\}$$

$$S_y(i+1) = \max_{v \in S} \{ S_v(i) \cdot P(y_{i+1}=y | Y_i=v) \cdot P(x_{i+1}=x_{i+1} | Y_{i+1}=y) \}$$

$$S_y(1) = P(Y_1=y) P(X_1=x_1 | Y_1=y)$$

### Trellis Graph



S 1 2 3 4

DET	0.003 (1.0 per)	$2.82 \times 10^{-5}$ (5.0 per)		
PRP	0.282 (1.0 per)	$2.82 \times 10^{-5}$		
N	0.001 (1.1 per)	$4.128 \times 10^{-3}$	$1.5 \times 10^{-3}$	
PREP	0.001 (1.1 per)	$5.64 \times 10^{-4}$	$1.25 \times 10^{-4}$	
V	0.002 (4.0 per)	$8.69 \times 10^{-2}$		
	PRP	V	Prep	N

$$\sum_{\text{DET}} (1) = 0.3 \times 0.01 = 0.003$$

$$\sum_{\text{PRP}} (1) = 0.3 \times 0.94 = 0.282$$

$$\sum_{\text{N}} (1) = 0.1 \times 0.01 = 0.001$$

$$\sum_{\text{PREP}} (1) = 0.1 \times 0.01 = 0.001$$

$$\sum_{\text{V}} (1) = 0.2 \times 0.01$$

$$\sum_{\text{DET}} (2) = \max \left\{ \begin{array}{l} V = \text{DET} \quad 0.003 \times 0.01 \times 0.01 \\ V = \text{PRP} \quad 0.282 \times 0.01 \times 0.01 \\ V = \text{N} \quad 0.001 \times 0.01 \times 0.01 \\ V = \text{PREP} \quad 0.001 \times 0.3 \times 0.01 \\ V = \text{Verb} \quad 0.002 \times 0.2 \times 0.01 \end{array} \right\} = 2.82 \times 10^{-5}$$

$$\sum_{\text{PRP}} (2) = \max \left\{ \begin{array}{l} \text{DET} \quad 0.003 \times 0.01 \times 0.01 \\ \text{PRP} \quad 0.282 \times 0.01 \times 0.01 \\ \text{N} \quad 0.001 \times 0.2 \times 0.01 \\ \text{PREP} \quad 0.001 \times 0.2 \times 0.01 \\ \text{V} \quad 0.002 \times 0.19 \times 0.01 \end{array} \right\} = 2.82 \times 10^{-5} \text{ PRP}$$

$$Y^* = (\text{PRP}, \text{V}, \text{PREP}, \text{N})$$

L#18

## STATISTICAL Learning Theory: PAC Learning

### Review:

- Structured Prediction  
→ sequence prediction  
 $h: X \rightarrow Y$  ( $X, Y$ ) are sequences

$Y$  = structured multivariate object

- HMM Generative Model  
 $P(X, Y) = P(y_1) P(x_1 | y_1) \prod_{i=2}^l P(y_i | y_{i-1}) P(x_i | y_i)$

- Compute  $h(x) = \arg \max_y P(x, y)$

- Viterbi Alg  $O(l | \Sigma_y |^2)$

$l$  = length of sentence  
 $|\Sigma_y|$  = # of tags

- Estimating  $P(X, Y)$  is extremely difficult b/c  $Y$  is sequence.  
therefore we have Markov assumption  $P(x_i | y_{i-1}) \leq P(x_i | y_i)$ . After

this assumption, we can easily compute  $\arg\max_y P(X, y)$  using Viterbi algorithm.

Setup:

- $|H|$  students  $H = \{h_1, \dots, h_{|H|}\}$
- $n$  bit sequence
- $p = 0.5$  probability of erring on a single bit given you are not psychic

Question: How likely is it that a particular player " $i$ " guesses the code without being psychic?

$$P(h_i \text{ correct} | h_i \text{ nonpsychic}) = (1-p)^n$$

$$\vdots$$
$$P(h_{|H|} \text{ correct} | h_{|H|} \text{ nonpsychic}) = (1-p)^n$$

Real Question: How likely is it that at least one student ~~guesses~~ guesses code without anybody being psychic?

$$P(h_1 \text{ correct} \vee h_2 \text{ correct} \vee \dots \vee h_{|H|} \text{ correct} | \text{all non psychic})$$

$$= 1 - P(h_1 \text{ not correct} \wedge \dots \wedge h_{|H|} \text{ not correct} | \text{all non psychic})$$

with independent assumption,

$$= 1 - \left[ P(h_1 \text{ not correct} | h_1 \text{ not psychic}) \dots P(h_{|H|} \text{ not correct} | h_{|H|} \text{ not psychic}) \right]$$

$$= 1 - \underbrace{\left[ (1 - (1-p)^n) \dots (1 - (1-p)^n) \right]}_{|H| \text{ times}}$$

$$= 1 - \prod_{i=1}^{|H|} [1 - (1-p)^n]$$

$$= 1 - (1 - (1-p)^n)^{|H|}$$

$$\text{let } |H| = 173, n = 4, p = \frac{1}{2}: 1 - \left[ 1 - \left( \frac{1}{2} \right)^4 \right]^{173} = 0.9999986$$

i.e. I am pretty sure that somebody will guess code. There is no evidence that all are non psychic.

How long do I have to make the sequence to be confident of psychic abilities?

$1 - (1 - (1-p)^n)^{|H|} < \delta$  small i.e. we want to find "n" so that prob of getting tracked goes down

$$n > \log_{(1-p)} \left[ 1 - (1-\delta)^{\frac{1}{|H|}} \right]$$

for  $\delta = 0.05$ , (prob of detecting n bits given every one is non-psychic)

$$|H| = 173$$

$$p = 0.5 \text{ (null hypothesis)}$$

$$n > \log_{1/2} (1 - 0.95^{1/173}) = 11.7$$

i.e. we have to make our sequence 12 bits long to make the chance of detecting 12 bit sequence 5%.

### ANALOGY In Learning

- $H$  hypothesis space  $\hat{=}$  students in class
- $n$  size of training sample  $\hat{=}$  length of sequences
- $\text{Err}_{\text{train}}(h)$  number of train error  $\hat{=}$  number of bit guessed correctly
- $\text{Err}_p(h)$  prediction/true error  $\hat{=}$  p(prop err on bit)

Question: Can we bound (upper bound) the probability that our learning alg "L" returns the hypothesis with large prediction error ( $\text{Err}_p(h) \geq \epsilon$ )?

- Assumption:
- Learning algorithm "L" has finite hypothesis space  $H$  [for instance decision trees]
  - At least one hypothesis  $h \in H$  has zero prediction error  $\text{Err}_p(h) = 0 \rightarrow \text{Err}_{\text{train}}(h) = 0$
  - Learning Alg "L" returns zero training error hypothesis  $\hat{h}$ .

To answer the above question, let's decompose that question

- (1) Probability that single hypothesis " $h$ " with  $\text{Err}_p(h) \geq \epsilon$  has train error 0?

$$P(\text{Err}_S(h) = 0 \mid \text{Err}_p(h) \geq \epsilon) \leq (1 - \epsilon)^n$$

(our assumption lies in i.i.d data)

- (2) Probability that there exists at least one  $h \in H$  with  $\text{Err}_p(h) \geq \epsilon$  and no training error? (It makes it indistinguishable from  $h \in H$  that had  $\text{Err}_S(h) = \text{Err}_p(h) = 0$ ).

Let  $h_1 \dots h_K$  be the hypothesis in  $H$  with  $\text{Err}_p(h_i) \geq \epsilon$

$$P(\text{Err}_S(h_1) = 0 \vee \dots \vee \text{Err}_S(h_K) = 0 \mid \text{Err}_p(h_1) \geq \epsilon, \dots, \text{Err}_p(h_K) \geq \epsilon) \\ \leq \sum_{i=1}^K P(\text{Err}_S(h_i) = 0 \mid \text{Err}_p(h_i) \geq \epsilon) : \text{Here we have}$$

used union bound that does not have any independence assumption i.e.  $P(X_1 = x_1 \vee X_2 = x_2 \vee \dots \vee X_n = x_n) \leq \sum_{i=1}^n P(X_i = x_i)$ . Here we don't have independence assumption. ~~Imagine~~ Imagine two decision trees that differ by leaf, they will have dependence among their predictions b/c of their similarity in structure.

$$\leq K(1 - \epsilon)^n$$

here  $K$  is the no of hypothesis in hypothesis space

$$K(1 - \epsilon)^n \leq |H|(1 - \epsilon)^n$$

How to upper bound  $K$ ? we can upper bound it by  $|H|$ .

$$\cancel{|H|(1 - \epsilon)^n} \quad K(1 - \epsilon)^n \leq |H|(1 - \epsilon)^n \leq |H|e^{-\epsilon n}$$

Hypothesis Space

$$\{h_i : \text{Err}_p(h_i) < \epsilon\} \leftarrow \text{good hyp}$$

$$\{h_i : \text{Err}_p(h_i) > \epsilon\} \leftarrow \text{bad hyp space}$$

$h^*$  is perfect hyp?

$$\{h \in H : \text{Err}_S(h) = 0\}$$

BASICALLY, what we have proved is the probability that learning error returns hypothesis is actually from the set where  $h : \text{Err}_p(h) < \epsilon$  i.e. set of those hypothesis whose prediction error  $\text{Err}_p(h)$  is less than  $\epsilon$ .

→ The probability that  $h$  returns  $h$  with prediction error greater than  $\epsilon$  is at most  $|H| e^{-\epsilon n}$

Here  $n \rightarrow$  # of training examples

$|H| \rightarrow$  cardinality of hypothesis

Now we have an upper bound on the worst performance of algorithm "A".

## Sample Complexity

How many examples does a learning algorithm "A" need so that with prob  $1-\delta$  it learns  $h$  with pred error  $[Err_p(h)]$  less than  $\epsilon$ ?

$$|H| e^{-\epsilon n} \leq \delta$$

Let  $\delta = 0.05$  (prob of screwing up)  
now solve for "n".

$$\boxed{n \geq \frac{1}{\epsilon} (\log |H| - \log \delta)}$$

(under three assumptions we stated before)

All of this is PAC learning = Probably approximately correctly.  
We want <sup>approximately</sup> correct & we can be correct upto some measure of  $\epsilon$ . And we can never be fully correct i.e. we can always be probably correct with probability of messing up  $\delta$ . This is PAC learning.

i.e.  $P(Err_p(h_{k(s)}) \leq \epsilon) \geq 1 - \delta$

i.e. probability of learning hypothesis that have prediction error less than  $\epsilon$

is  $1 -$  prob of messing up.

By increasing "n", we can ~~also~~ increase the chances of getting such hypothesis whose

$$P(Err_p(h_{k(s)}) \leq \epsilon) \geq 1 - \delta.$$

This is PAC-learning.

# L#19: Statistical Learning Theory: Error bounds & VC-Dimension

Review:

- PAC ~~style~~ learning

$$P(\text{Err}_p(h_{L(S)}) \leq \epsilon) \geq 1 - \delta$$

-  $\epsilon$  is the amount of error (prediction error) we are willing to tolerate

-  $\delta$  is the prob that " $L$ " fail. It is the prob we allow learner to fail.

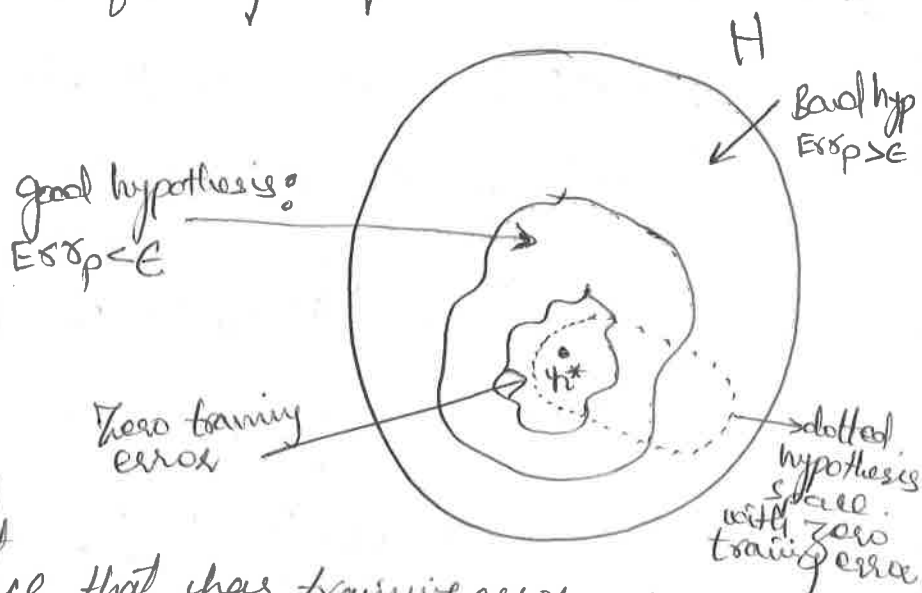
- We can never be sure about zero-generalization error.

- Generalization error bound for finite  $H$ , ~~the~~ concept is in  $H$ :

$$P(\text{Err}_p(h_{L(S)}) \geq \epsilon) \leq |H| \cdot e^{-\epsilon n}$$

$n = \#$  of training examples

\* Problem starts when our dotted hypothesis space actually goes out from the region of good hypothesis space where  $\text{Err}_p < \epsilon$ . Therefore, our learner " $L$ " can actually return hypothesis " $h$ " that lies in bad hypothesis space that has training error = 0 but  $\text{Err}_p > \epsilon$ .



The proof we did last time said:

"Can we ~~say~~ with certain probability that zero training error hypothesis are strictly contained in a set of good hypothesis set."

$$P(\text{Err}_p(h_{L(S)}) \leq \epsilon) \leq |H| e^{-\epsilon n}$$

"probability that no bad hypothesis ended up in the set of zero training error hypothesis".

## - Sample Complexity:

"How many training examples do I need to be reasonably sure that I get good hyp"

or

"It asks how many training examples do I need so that with probability  $1-\delta$  I get a hypothesis returned by my learning algorithm that has prediction error no greater than  $\epsilon$ ." Basically it solves for "n" the following

$$P[\text{Err}_S(h_{\text{LCS}}) \geq \epsilon] \leq |H| e^{-\epsilon n}$$

In today's lecture we will relax some assumptions:

- i) we can have hypothesis with non-zero training error
- ii) our hypothesis space is infinite i.e linear function.

Generalization Error bound: finite H, non-zero error

learner:  $L$  returns  $h \in H$  with minimum training error

Idea: Bound the deviation b/w training and generalization error

$$P(|\text{Err}_S(h_{\text{LCS}}) - \text{Err}_P(h_{\text{LCS}})| \leq \epsilon) \geq 1-\delta$$

\* we will say something about all hypotheses in hypothesis space

$$\Leftarrow P(\underbrace{\max_{h \in H} |\text{Err}_S(h) - \text{Err}_P(h)|}_{\text{what is the maximum deviation}} \leq \epsilon) \geq 1-\delta$$

$$\Leftrightarrow P(\max_{h \in H} |\text{Err}_S(h) - \text{Err}_P(h)| > \epsilon) < \delta \quad (\text{taking complement of above})$$

$$P(|\text{Err}_S(h_1) - \text{Err}_P(h_1)| > \epsilon \vee \dots \vee |\text{Err}_S(h_{|H|}) - \text{Err}_P(h_{|H|})| > \epsilon)$$

Using Union bound

$$\leq \sum_{i=1}^{|H|} P(|\text{Err}_S(h_i) - \text{Err}_P(h_i)| > \epsilon) \quad : \text{Err}_S(h_1) = \frac{1}{n} \sum_j \Delta(h(x_j), y_j)$$



Using Hoeffding/Chernoff Bound (for binomial distribution)

$$P\left(\left|\frac{1}{n} \sum_{i=1}^n x_i - p\right| \geq \epsilon\right) \leq 2e^{-2n\epsilon^2}$$

$$\leq \sum_{i=1}^{|H|} 2e^{-2n\epsilon^2} = |H| 2e^{-2n\epsilon^2}$$

We can turn this ~~into~~ around in terms of sample complexity.

Sample Complexity

$$|H| 2e^{-2n\epsilon^2} \leq \delta$$

$$\text{i.e. } n \geq \frac{1}{2\epsilon^2} (\log |H| - \log \delta - \log 2)$$

Example: Boolean conjunction over  $\ell$  literals ( $N=10, \ell=3$ )

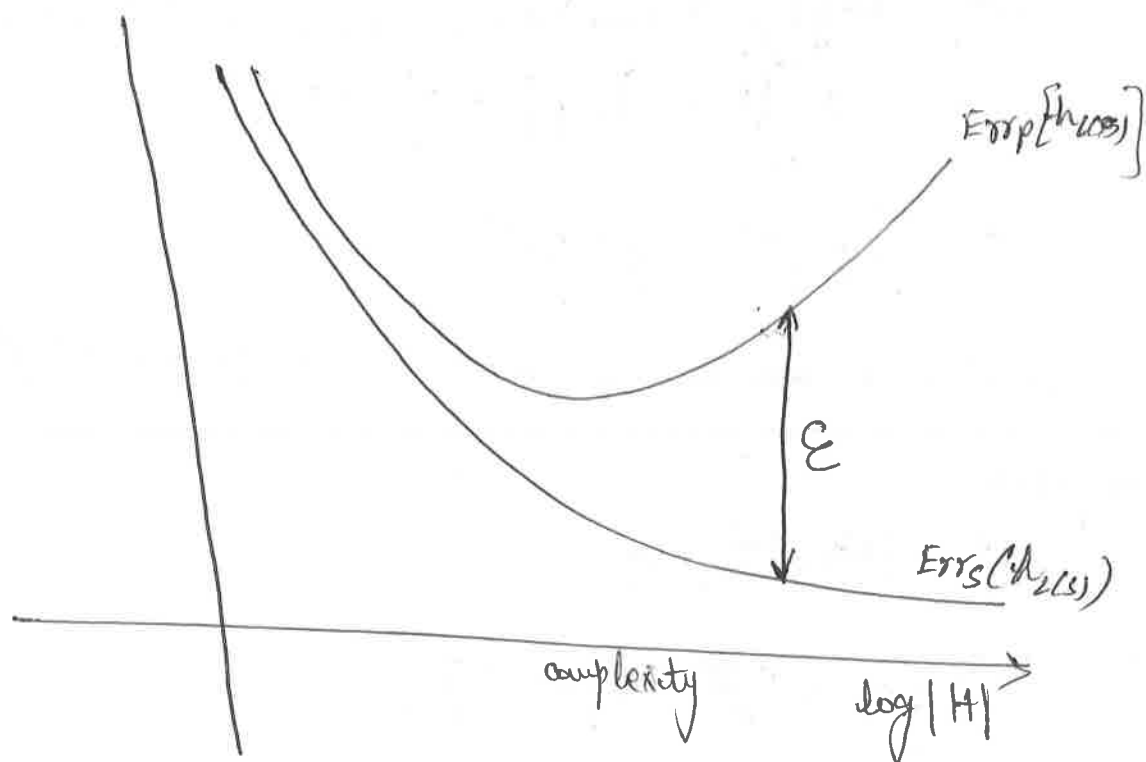
$$\rightarrow |H| = \binom{N}{\ell} \cdot 2^\ell = 960 \text{ Hypotheses}$$

How many examples do I need so that  $\epsilon = 0.1$  i.e. 10% deviation b/w training & prediction error?

$$\epsilon = 0.1$$

$$\delta = 0.05 \text{ (screw-up probability)}$$

$\Rightarrow n \geq 528$  (I need 529 examples so that I get no more than 10% deviation b/w training & prediction error with 95% confidence interval). In theory, it tells us how many examples do we need. In practice, we need much more example. These bound tells you behaviour.



These bound tells us why overfitting is happening.

Overfitting

for a fixed  $\delta$ , solve for  $\epsilon$

$$2|H|2e^{-2n\epsilon^2} = \delta$$

$$\Leftrightarrow \epsilon = \sqrt{\frac{1}{2n} \log(2|H|) - \log(\delta)}$$

$$\Rightarrow P(|\text{Err}_S(h_{L(S)}) - \text{Err}_P(h_{L(S)})| \geq \sqrt{\frac{1}{2n} (\log(2|H|) - \log(\delta))}) \leq \delta$$

Take complement

u

u

$\leq$

u.

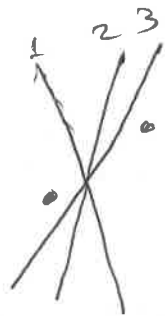
$) \geq (1-\delta)$

$$\Leftrightarrow \text{with prob } 1-\delta: \text{Err}_P(h_{L(S)}) \leq \text{Err}_S(h_{L(S)}) + \sqrt{\frac{1}{2n} (\log(2|H|) - \log \delta)}$$

## Infinite Hypothesis Spaces

Idea: Consider only the "effective" number of Hypothesis (i.e. hypothesis that don't give the same classification).

Example:



• we will lump 1, 2, 3 together b/c all of them classify exactly similarly

Question: How many different hypothesis exist in  $H$  for sample  $S$ ?

$\Pi_H(S)$  contains one  $h \in H$  for each different classification of sample " $S$ "

→ <sup>It is</sup> effectively finite number of hypothesis

(• If we have " $n$ " points, we will have  $2^n$  different hypothesis)

→ fortunately, for many infinite  $H$   
 $|\Pi_H(S)| \ll 2^n$  for all samples  $S$ .

L#20

## STATISTICAL Learning theory: Experts And Bandit

Review:

- Generalization Error Bound

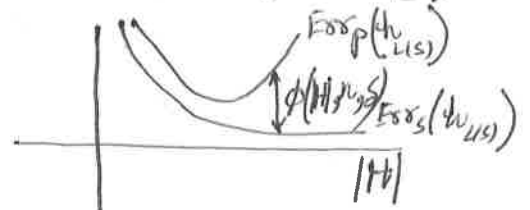
→ bound  $P(|Err_S(h_{L(S)}) - Err_P(h_{L(S)})| \geq \epsilon) \leq \delta$

→ with prob  $1 - \delta$ ,  $Err_P(h_{L(S)}) \leq Err_S(h_{L(S)}) + \phi(|H|, n, \delta)$

- Infinite Hypothesis Spaces

$\Pi_H(S)$  contains one  $h \in H$  for each different classification

of  $S$ . It is always a finite number and usually  $\ll 2^n$



Consider this datapoint:



we have 5 points. We need linear hyperplane to separate them. Thus  $\Pi_H(S)$  is much less than  $2^n$ .

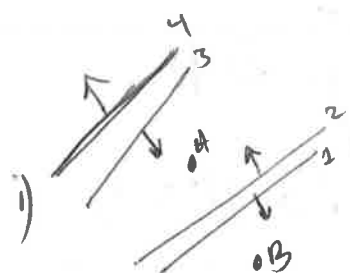
Sauer's Lemma

$$\forall S, |\Pi_H(S)| \leq \left( \frac{e n}{VCdim(H)} \right)^{VCdim(H)}$$

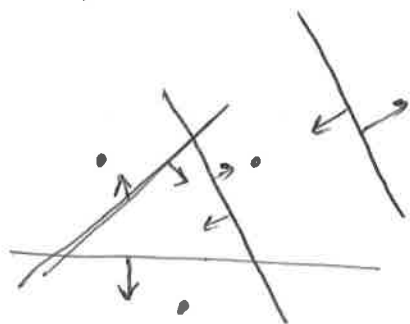
$|\Pi_H(S)|$  is effective measure of hypo.

Example: VC dimension of Hyperplanes in  $\mathbb{R}^2$

$$H_{L2} = \{h : h(x) = \text{sign}(\omega \cdot x + b)\}$$



$$VC\text{-dimension}(H_{L2}) \geq 2$$



shattered three points with 3 hyperplanes

$$VCdim(H_{L2}) \geq 3$$



we cannot ~~class~~ separate B, C from A & D using linear hyperplane



i.e you cannot shatter 4 points in 2D planes.

$$\boxed{VC \dim(H_{22}) < 4}$$

General:

- For linear classifiers in  $\mathbb{R}^N$

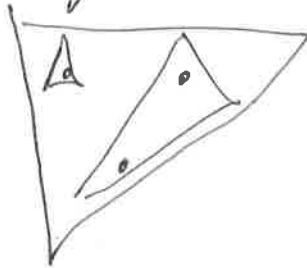
$$VC \dim(H_{LN}) = N+1$$

- For linear classifiers with margin  $\gamma$

$$VC \dim(H_\gamma) \leq \frac{R^2}{\gamma^2}$$

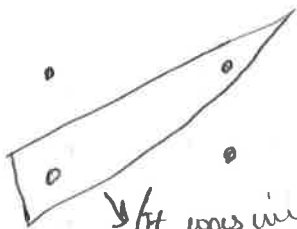
VC dimension of Triangles in  $\mathbb{R}^2$

i)



Can a triangle shatters three points?  
A: Yes

ii)

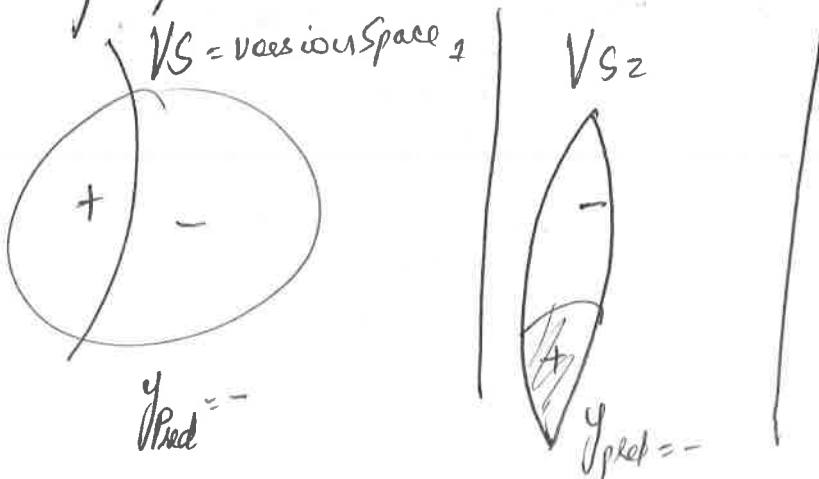


Can a triangle shatters 4 points?  
A: Yes

(It was impossible in linear classifiers)

And move on to find 5, 6 ... points to find the lower bound of VC dimension of triangles in  $\mathbb{R}^2$ .

Halving Algorithm



Theorem: After  $\log_2(N)$  mistakes, the version space contains only the perfect expert.

Proof: Consider each iteration  $t$ :

- Case 1: Algorithm predicts incorrectly  
→ more than half of the hypothesis in  $VS_t$  were wrong & kicked out.

$$\rightarrow |VS_{t+1}| \leq \frac{1}{2} |VS_t|$$

- Case 2: Algorithm predicts correctly

$$\rightarrow |VS_{t+1}| \leq |VS_t|$$

→ Can only ~~half~~ halve VS for  $\log_2(N)$  times before we run out of hypothesis.

→ It is with no noise. Can we make it robust to noise? See Weighted Majority Algorithm.

Example:  $H = \{h_1, h_2, h_3\}$ ,  $\beta = \frac{1}{2}$  (expects = hypothesis)

-  $t=1$ : given  $x_1$ :  $h_1(x_1)=1$ ,  $h_2(x_1)=1$ ,  $h_3(x_1)=-1$

$w_1 = (1, 1, 1)$ : Weight of  $y=+1$  is  $(w_{11}+w_{12})=2$   
weight of  $y=-1$  is  $w_{13}=1$  } Prediction  $y=+1$

Receive true label  $y_1 = y_1 = -1$  (we made an error)  
 $w_2 = (\frac{1}{2}, \frac{1}{2}, 1)$

-  $t=2$ : given  $x_2$ :  $h_1(x_2)=-1$ ,  $h_2(x_2)=1$ ,  $h_3(x_2)=-1$

$w_2 = (\frac{1}{2}, \frac{1}{2}, 1)$ : weight of  $y=+1$  is  $\frac{1}{2}$   
weight of  $y=-1$  is  $\frac{3}{2}$  }  $y_{\text{pred}} = -1$

Receive  $y_{\text{true}} = -1$

$$w = (\frac{1}{2}, \frac{1}{4}, 2)$$

i.e. we can see hypothesis that predicts correctly have larger weight.

Theorem: Let  $\Delta^*$  be the number of mistakes of the best expert in hindsight (experience), the weighted majority algorithm makes at most  $\Delta^{WM} \leq 2.4 (\Delta^* + \log_2(N))$  for  $\beta = \frac{1}{2}$

Proof: Let  $h^*$  be the best expert in hindsight  $\rightarrow w^* = (\frac{1}{2})^{\Delta^*}$

Let  $W_t = \prod_{i=1}^t W_{t+1}$ .  $W_1 = n$ . for each mistake,  $W_{t+1} \leq \frac{3}{4} W_t$ , since more than  $\frac{1}{2}$  half of the weight gets slashed by  $\beta = \frac{1}{2}$ .

Since  $h^* \in H$ :  $\underbrace{(\frac{1}{2})^{\Delta^*}}_{\text{one expert}} \leq \underbrace{n (\frac{3}{4})^{\Delta^{WM}}}_{\text{sum of all other experts}} \rightarrow \text{solve for } \Delta^{WM} \leq \frac{\Delta^* + \log_2(N)}{-\log_2(\frac{3}{4})}$

L#21

## Clustering: Similarity based Clustering

### Review

- VC dim of  $H$ 
  - largest set of examples that can be shattered by functions from hypothesis space  $H$ .
- Sour's lemma
  - Bound:  $Err_P(h_{crs}) \leq Err_S(h_{crs}) + \phi\left(\sqrt{\frac{VC \cdot \dim(H) + \log \epsilon}{n}}\right)$
- Expert Setting
  - Halving Algorithm (perfect expert exist, 0/1 loss)  $\rightarrow \log_2(|H|)$  mistakes
  - Weighted majority algorithm (0/1 loss)  $\rightarrow \Delta^{WM} \leq 2.4 (\Delta^* + \log_2(|H|))$   
 $\Delta^* = \text{loss of experts in hindsight}$
  - Exponentiated Gradient Algorithm (bounded loss)
    - $\Delta^{EG} \leq \Delta^* + O(\sqrt{n \log(|H|)})$
    - Divide by "n" to get upper bound on VC-dimension
    - $\frac{\Delta^{EG}}{n} \leq \frac{\Delta^*}{n} + O\left(\sqrt{\frac{\log(|H|)}{n}}\right)$

- Bandit Setting
- EXP3 algorithm

## Supervised learning

- $P(X, Y)$  learning task
- i.i.d training sample  $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$
- learn  $h: X \rightarrow Y$

## Clustering: (trying to find structure in $P(X)$ )

- $P(X)$
- i.i.d sample  $S = \{x_1, \dots, x_n\}$
- To learn: structure of  $P(X)$

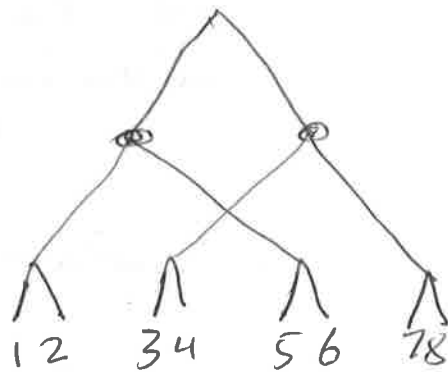
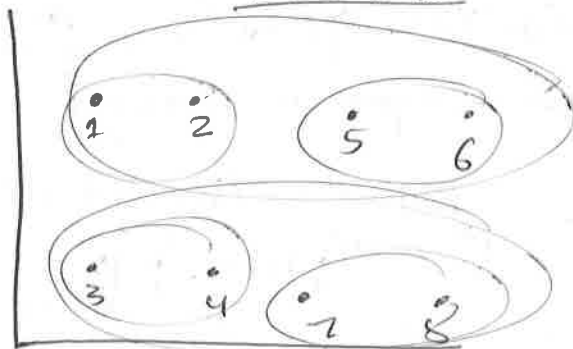
## Outliers Detection:

- $P(X)$ : we don't know
- i.i.d sample  $S = (x_1, \dots, x_n, x'_1, \dots, x'_m)$
- find  $x'$  that did not come from  $P(X)$

## Novelty Detection:

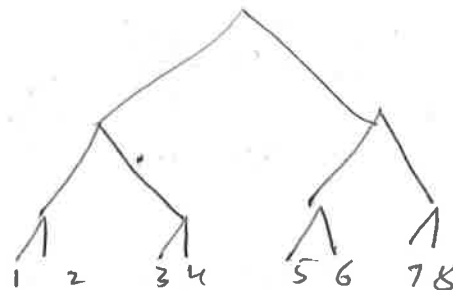
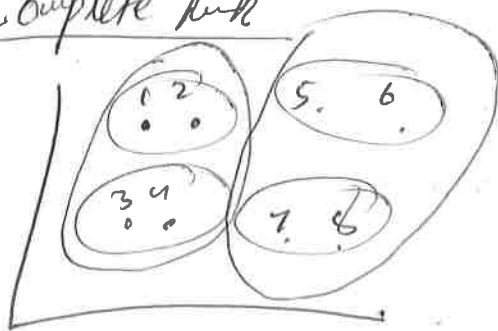
- $P(X)$  and  $P(X')$
- $S = (x_1, \dots, x_k, x'_1, \dots, x'_m)$   
we want to find where the switch happens
- find locations of change

close-link





## Complete link



## L#22: Clustering: K-Means & Mixture of Gaussians

### Review:

- Unsupervised learning  $\rightarrow$  Clustering (finding structure in distributions)
- Hierarchical Aggl Clustering
  - $\rightarrow$  start with instances, then merge
  - $\rightarrow$  merging gives dendrogram (tree)
- Similarity among clusters
  - $\rightarrow$  single link (similarity b/w two most similar members of clusters)
  - $\rightarrow$  complete-link ( " " " least " " " "
  - $\rightarrow$  group average
- $O(n^2 \log n)$  for single/complete link

Non-Hierarchical (flat) Clustering  $\rightarrow$  K-Means  
 $\rightarrow$  GMM with EM

### Objective function of K-Means

$$\min_{C=\{C_1, \dots, C_K\}} \left\{ \sum_{j=1}^K \sum_{x_i \in C_j} \|x_i - \mu_j\|^2 \right\}$$

It is redundant  $\rightarrow \mu_1, \dots, \mu_K$

Step 1: update centroids

$$\text{Given } C=\{C_1, \dots, C_K\}$$

$$\left\{ \min_{\mu_j} \sum_{x_i \in C_j} \|x_i - \mu_j\|^2 \right\}$$

$\rightarrow$  solution of this is centroid. You can take derivative & set to 0

$$\min_{x \in \mu_j} \sum_{x_i \in C_j} \|x_i - \mu_j\|^2 = \frac{1}{|C_j|} \sum_{x_i \in C_j} x_i$$

→ Here K-Means improves objectives

Step 2: Assign points to closest mean

→ K-Means improves objective

In both steps, K-Means improves objective function

Theorem: K-Means always terminates

Proof: finite number of clusters and no cycles since objective always improves.

### Clustering as Prediction

→  $P(X)$  unknown

→  $S = (x_1, \dots, x_n)$

→  $H = \{h_1, \dots, h_n\}$  where  $h_i$  is different clustering of  $S$ .

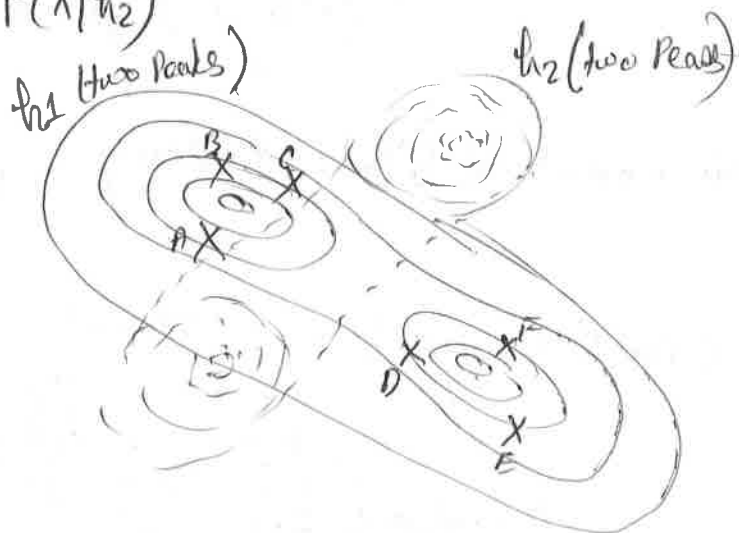
We need to find which  $h_i$  in  $H$  generated our data  $S$ ?

Example: Pick b/w two hypothesis:

$P(X|h_1)$  and  $P(X|h_2)$

i) what's the prob that  $h_1$  generated A, B, C, D, E, F?  
Ans: Very high

ii) what's the probability that  $h_2$  generated A...F  
Ans: Very low



Why is  $h_1$  better than  $h_2$ ?

$$P(S|h_1) = \prod_{i=1}^n P(x_i|h_1) \quad \text{pretty high}$$

$$P(S|h_2) = \prod_{i=1}^n P(x_i|h_2) \quad \text{very low}$$

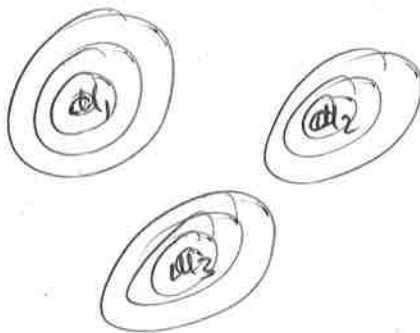
$P(S|h) =$  likelihood function

Strategy

we will pick  $h_i$  with maximum likelihood. Maximum likelihood inference is very widespread in statistics.

Generative Process of Mixture of Gaussians

lets assume mixture of 3 Gaussian distributions.



~~Gauss~~ GMM will be parametrized by  $(\mu_1, \mu_2, \mu_3, \Sigma_1, \Sigma_2, \Sigma_3)$

we will toss a dice (three sides), then we with probability  $P(Y=j)$ . Then we will select that particular mixture and sample point  $x_i$  from that Gaussian  $N(x_i; \mu_j, \Sigma_j)$ .

GMM

$$\theta = (\mu_1, \mu_2, \dots, \mu_K, \Sigma_1, \dots, \Sigma_K)$$

has to be estimated by EM

\* EM is used for latent variable models.

GOAL: Find  $\mu_1, \dots, \mu_K$  for which sample "S" has maximum likelihood.

$$\max_{\mu_1, \dots, \mu_K} P(S | \mu_1, \dots, \mu_K) = \max_{\mu_1, \dots, \mu_K} \prod_{i=1}^n \sum_{j=1}^K P(X=x_i | y=j, \mu_j) P(y=j)$$

$$\frac{d}{d\mu_j} (\text{obj}) = 0$$

$$\forall j \frac{d}{d\mu_j} \left[ \prod_{i=1}^n \sum_{l=1}^K P(X=x_i | y=l, \mu_l) P(y=l) \right] = 0$$

$$\Leftrightarrow \forall j \frac{d}{d\mu_j} \left[ \sum_{i=1}^n \frac{d}{d\mu_j} \left[ \ln \left( \sum_{l=1}^K P(X=x_i | y=l, \mu_l) P(y=l) \right) \right] \right] = 0$$

$$\Leftrightarrow \forall j \sum_{i=1}^n \frac{1}{\sum_{l=1}^K P(X=x_i | y=l, \mu_l) P(y=l)} \frac{d}{d\mu_j} \left[ \sum_{l=1}^K P(X=x_i | y=l, \mu_l) P(y=l) \right]$$

$$\Leftrightarrow \forall j \sum_{i=1}^n \frac{1}{\dots} \frac{d}{d\mu_j} \left[ P(X=x_i | y=l, \mu_l) P(y=l) \right]$$

$$\Leftrightarrow \forall j \sum_{i=1}^n \frac{P(X=x_i | y=l, \mu_l) P(y=l)}{\dots} \times \frac{d}{d\mu_j} \left[ \frac{1}{2} (x_i - \mu_j)^2 \right] \left[ P(X=x | \mu) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-\mu)^2} \right]$$

$$\Leftrightarrow \forall \frac{d}{d\mu_j} P(y=l | X=x_i, \mu_j) (x_i - \mu_j) = 0$$

→ solve for  $\mu_j$

$$\mu_j = \dots \text{ in slide}$$

$\mu_j$  = weighted average : we are weighting the points by probabilities. In K-Means, we were measuring distance

\* if we assign probability  $P(Y=j | X=x_i, \mu) = 1$  to the cluster it ~~assigns~~ belongs, we can turn EM to K-Means. In EM, we have

soft assignment of clusters rather than hard assignment of K-means.

## L#23: structured Output Prediction

\* Can we learn HMM in discriminative fashion?

$$P(x, y) = P(y_1) \cdot P(x_1 | y_1) \cdot \prod_{i=2}^l P(y_i | y_{i-1}) P(x_i | y_i)$$

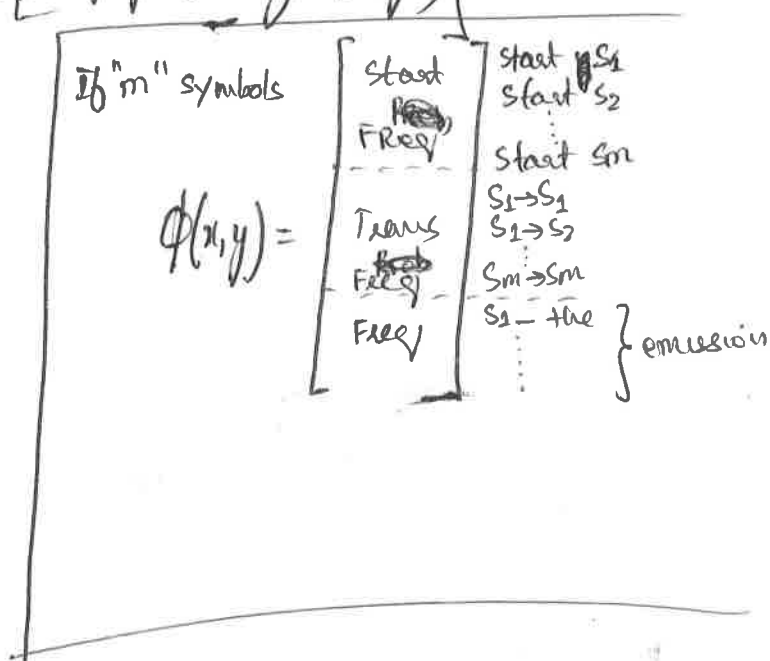
Apply log

$$\log P(x, y) = \log P(y_1) + \log P(x_1 | y_1) + \sum_{i=2}^l [\log P(y_i | y_{i-1}) + \log P(x_i | y_i)]$$

$$= w^T \phi(x, y)$$

"w" will be the log probabilities

In this way  ~~$\arg\max_y P(x, y)$~~   $\arg\max_y w^T \phi(x, y)$   
 $= \arg\max_y P(y | x)$



\* Example

the/det bank/N opens/V

$$\phi(x, y) = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} \begin{matrix} \text{DET} \\ \text{DET} \rightarrow N \\ N \rightarrow V \\ \text{DET} - the \\ N - bank \\ V - opens \end{matrix}$$

• W

the/det bank/V opens/V

$$\phi(x, y) = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} \begin{matrix} \text{det} \\ \text{det} \rightarrow N \\ V \rightarrow V \\ \text{det} - the \\ \text{bank} \rightarrow \text{bank} \\ V \rightarrow \text{opens} \end{matrix}$$

Different than above  $\phi(x, y)$

Thus we get linear discriminative function and we will minimize error on training set.