

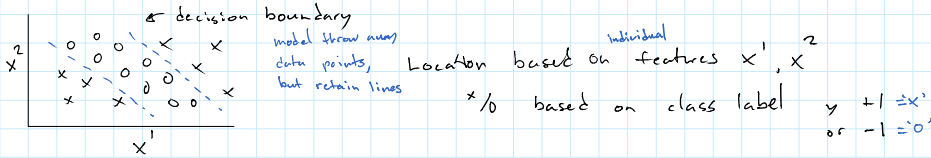
Supervised Learning

Input: X, Y design matrix, target

Output function: $f(x_i) \mapsto y_i$

Evaluation compare: $f(x_i) = \hat{y}_i$ to (real) y_i on new (x_i, y_i)

Notation: x_i is training example 'i' (vector) features
 x^i is variable 'i' (scalar) feature



KNN 'closest' - euclidean distance $d(\vec{w}, \vec{v}) = \sqrt{\sum_{i=1}^n (w_i - v_i)^2} = \|\mathbf{w} - \mathbf{v}\|$
 - others possible (?) L1 norm?

KNN: 'cost'

- checking $y_i - x_j$ is $O(D)$ (?.D) = distance
- do this N times for each of T test params $O(D \cdot N \cdot T)$
- sorting $O(N \log N)$
- selecting $O(KN)$

Possible course project: Improve KNN
 - Better distance define or learn distance
 - faster classification - gpu, tree, sparse

Q: why can we learn?



- Probability

assume $(x_i, y_i) \sim D$ sampled by same distribution, we learn the distribution

$0 \leq P(A) \leq 1$ definitely not, definitely

$P(\neg A) = 1 - P(A)$

$P(x) \stackrel{\Delta}{=} P(X=x)$

$\sum_{x \in X} P(x) = 1$

$P(A, B) = P(A \cap B)$

$P(A \cup B) = P(A) + P(B) + P(A, B)$

$P(A|B) = \frac{P(A, B)}{P(B)}$ (for $P(B) > 0$)

$P(A, B) = P(A|B)P(B) = P(B|A)P(A)$

Bayes Rule

$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$

$\propto P(B|A)P(A)$ $\propto P(B)$ is fixed

$P(A) = \sum_b P(A, B=b)$

$p(\text{even}) = \sum_{i=1}^6 p(\text{even}, i) = 0 + \frac{1}{6} + \dots = \frac{1}{2}$

We can always add a conditioning case.

$P(A|B, c) = \frac{P(B|A, c)P(A|c)}{P(B|c)}$

$X \perp Y \Leftrightarrow P(x, y) = P(x)P(y)$

because sum across $P(i) = 1$

$$\sum_{i=1}^6 \prod_{j=1}^6 P(x_j) P(i) = \frac{\prod_{j=1}^6 P(x_j)}{\prod_{j=1}^6 P(x_j)}$$

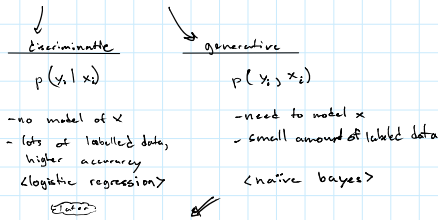
1	2	3	4
5	6	7	8
9	10	11	12

$p(1) = 1/6$
 $p(3) = 1/6$

$P(2, \text{even}) = 0$
 $P(3, \text{odd}) = 1/6$
 $P(1 \cup 2) = p(1) + p(2) - p(1, 2) = \frac{1}{6} + \frac{1}{6} - 0 = \frac{1}{3}$ (intuitive)

$P(3 | \text{odd}) = \frac{P(3, \text{odd})}{P(\text{odd})} = \frac{1/6}{1/2} = \frac{1}{3}$ (also intuitive)

Probabilistic Classifiers



Training
 $p(y_i, x_i) = p(x_i | y_i) p(y_i)$

Testing
 $p(y_i | x_i) = \frac{p(x_i | y_i) p(y_i)}{p(x_i)} = \propto p(x_i | y_i) p(y_i)$

↑ easy times we say +, rest times we say -

x_i^1	x_i^2	y_i
0	0	0
0	0	1
0	1	0
⋮		

↓ a lot, 2^{D+1} parameters

Naive Bayes

x^i are "mutually independent"

↑ subsets $X^A | B$ | y each x event is separate

$$P(x_i^1 | y_i) = P(x_i^1 | x_i^2=0, y_i) P(x_i^2=0 | y_i)$$

$$= P(x_i^1 | y_i) P(x_i^2=0 | y_i)$$

$$= \prod_{j=1}^D P(x_i^j | y_i)$$

x^j	y	Parametric model.
0	0	4D as opposed to
0	1	2^{D+1}
1	0	
1	1	

Maximum Likelihood

$$\text{argmax}_{\theta} p(D | \theta)$$

$$\Downarrow$$

$$\text{argmax}_{\theta} \log(D | \theta)$$

↑ how we parameterize our data

$$\Rightarrow \text{argmax}_{\theta} \prod_{i=1}^D p(\theta_i | \theta)$$

$$\Rightarrow \text{argmax}_{\theta} \sum_{i=1}^D \log p(\theta_i | \theta)$$

$$\log p(y_i, x_i) = \sum_{i=1}^n \log p(y_i | \theta) + \sum_{j=1}^D \log p(x_i^j | y_i, \theta_j^x)$$

$$(y_i | \theta) \sim \text{Bernoulli}(\theta) \quad [0, 1]$$

$$p(y_i | \theta) = \underbrace{\theta^{\sum y_i=1}}_t \underbrace{(1-\theta)^{\sum y_i=0}}_f \quad \text{one or the other}$$

$$\text{argmax}_{\theta} \sum_{i=1}^N \log p(y_i | \theta)$$

$$= \sum_{i=1}^N \mathbb{I}(y_i=1) \log \theta + \mathbb{I}(y_i=0) \log(1-\theta)$$

Derivative (gradient)

$$\nabla = \frac{N_1}{\theta} - \frac{N_0}{1-\theta}$$

find $\nabla = 0$:

$$(1-\theta) N_1 = N_0 (\theta)$$

$$\frac{\theta}{1-\theta} = \frac{N_1/N}{N_0/N}$$

$$N_1 + N_0 = N$$

$$\theta = \frac{N_1}{N}$$

$y_i \in \{1, 2, 3, 4\}$

$$\theta^c = \frac{N_c}{N}$$

$$(x_i | y_i, \theta_j^x)$$

$$\theta_j^c = \frac{N_{jc}}{N_c} \quad \text{feature } j \text{ and class } c$$