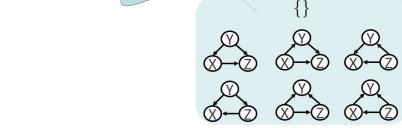
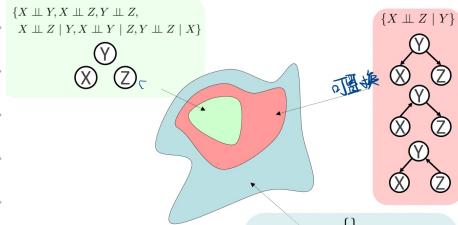


Possibilities

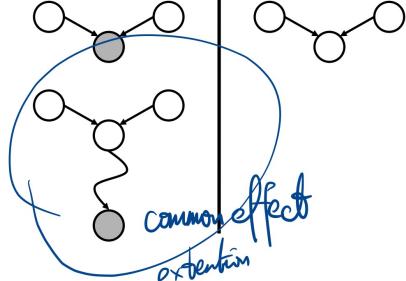
$$\text{Bayes Rule: } P(x|y) = \frac{P(y|x)}{P(y)} P(x)$$



Active Triples



Inactive Triples



IN: evidence instantiation

$$w = 1.0$$

for $i=1, 2, \dots, n$

if X_i is an evidence variable

x_i = observation for X_i

$$\text{Set } w = w * P(x_i | \text{Parents}(X_i))$$

else

$$\text{Sample } x_i \text{ from } P(X_i | \text{Parents}(X_i))$$

return $(x_1, x_2, \dots, x_n), w$

Markov Chain

$$P(X_t) = \sum_{X_{t-1}} P(X_{t-1}=x_{t-1}) P(X_t | X_{t-1}=x_{t-1})$$

$$P(X_{t+1}) = T^T P(X_t) \quad \downarrow \text{filtering}$$

$$P(X_{t+1} | e_{1:t+1}) = \alpha P(e_{t+1} | X_{t+1}) \sum_{x_t} P(x_t | e_{1:t}) P(X_{t+1} | x_t)$$

$$f_{1:t+1} = \alpha O_{t+1} T^T f_{1:t}$$

$$\text{Prediction: } P(x_{t+1} | e_{1:t}) = \sum_x P(x_{t+1} | x_t) P(x_t | e_{1:t})$$

Filtering: $P(X_t | e_{1:t})$

belief state — input to the decision process of a rational agent

Prediction: $P(X_{t+k} | e_{1:t})$ for $k > 0$

evaluation of possible action sequences; like filtering without the evidence

Smoothing: $P(x_k | e_{1:t})$ for $0 \leq k < t$

better estimate of past states, essential for learning

Most likely explanation: $\arg \max_{x_{1:t}} P(x_{1:t} | e_{1:t})$

speech recognition, decoding with a noisy channel

$$A \Rightarrow B \equiv \neg A \vee B$$

State Trellis

$$\text{value of arc: } P(x_t | x_{t-1}) P(e_t | x_t)$$

$$m_{1:t+1} = \text{VITERBI}(m_{1:t}, e_{t+1}) \\ = P(e_{t+1} | X_{t+1}) \max_{x_t} P(X_{t+1} | x_t) m_{1:t}$$

$$A: a_{ji} = P(X_{t+1} = j | X_t = i)$$

$X_t X_{t+1}$	A	B	H	S
A	0.6	0.1	0.1	0.2
B	0.0	0.3	0.2	0.5
H	0.8	0.1	0.1	0.1
S	0.2	0.0	0.1	0.7

$$B: b_{ik} = P(O_t = k | X_t = i)$$

$X_t O_t$	p	e	b	i
A	0.6	0.2	0.1	0.1
B	0.1	0.4	0.1	0.4
H	0.0	0.0	0.7	0.3
S	0.0	0.0	0.1	0.9

$$\pi = P(X_t = i):$$

A	B	H	S
0.5	0.0	0.0	0.5

Step 1: Initialize $\delta_1(i)$

$$O_1 = \underline{\underline{b}}$$

$$\delta_1(i) = \begin{bmatrix} 0.5 \\ 0.0 \\ 0.0 \\ 0.5 \end{bmatrix} \odot \begin{bmatrix} 0.1 \\ 0.7 \\ 0.1 \\ 0.5 \end{bmatrix} = \begin{bmatrix} 0.05 \\ 0.0 \\ 0.0 \\ 0.05 \end{bmatrix}$$

$$\pi_1 = \underline{\underline{0(e)}}$$

Observations:

$$o_{1:t} = \{b, p, l, e\}$$

Find:

Most likely hidden state sequence: $X_{1:t+1}^*$

$$O_2 = \underline{\underline{p}}$$

$$\delta_2(i) = \underline{\underline{ }}$$

$$\begin{aligned} A \rightarrow A & (p) & B \rightarrow A & H \rightarrow A & S \rightarrow A \\ \max(0.05 \times 0.6 \times 0.6, 0.0 \times 0.2 \times 0.6, 0.8 \times 0.8 \times 0.6, 0.05 \times 0.2 \times 0.6) & = & \max(0.05 \times 0.1 \times 0.1, 0.0 \times 0.4 \times 0.1, 0.0 \times 0.1 \times 0.1, 0.05 \times 0.0 \times 0.1) & = & \max(0.05 \times 0.1 \times 0.9, 0.0 \times 0.7 \times 0.9, 0.0 \times 0.1 \times 0.9, 0.05 \times 0.7 \times 0.9) & = \\ \max(0.05 \times 0.1 \times 0.1, 0.0 \times 0.4 \times 0.1, 0.0 \times 0.1 \times 0.1, 0.05 \times 0.0 \times 0.1) & = & \max(0.05 \times 0.1 \times 0.1, 0.0 \times 0.4 \times 0.1, 0.0 \times 0.1 \times 0.1, 0.05 \times 0.0 \times 0.1) & = & \max(0.05 \times 0.1 \times 0.9, 0.0 \times 0.7 \times 0.9, 0.0 \times 0.1 \times 0.9, 0.05 \times 0.7 \times 0.9) & = \\ \max(0.005 \times 0.1 \times 0.1, 0.0 \times 0.4 \times 0.1, 0.0 \times 0.1 \times 0.1, 0.0005 \times 0.0 \times 0.1) & = & \max(0.005 \times 0.1 \times 0.1, 0.0 \times 0.4 \times 0.1, 0.0 \times 0.1 \times 0.1, 0.0005 \times 0.0 \times 0.1) & = & \max(0.005 \times 0.1 \times 0.9, 0.0 \times 0.7 \times 0.9, 0.0 \times 0.1 \times 0.9, 0.0005 \times 0.7 \times 0.9) & = \\ \max(0.0005 \times 0.1 \times 0.1, 0.0 \times 0.4 \times 0.1, 0.0 \times 0.1 \times 0.1, 0.00005 \times 0.0 \times 0.1) & = & \max(0.0005 \times 0.1 \times 0.1, 0.0 \times 0.4 \times 0.1, 0.0 \times 0.1 \times 0.1, 0.00005 \times 0.0 \times 0.1) & = & \max(0.0005 \times 0.1 \times 0.9, 0.0 \times 0.7 \times 0.9, 0.0 \times 0.1 \times 0.9, 0.00005 \times 0.7 \times 0.9) & = \\ 0.0005 & = & 0.00005 & = & 0.000005 & = \end{aligned}$$

$$\delta_3(i) = \underline{\underline{l}}$$

$$\pi_2 = \underline{\underline{0(l)}}$$

$$\begin{aligned} \max(0.018 \times 0.6 \times 0.1, 0.0005 \times 0.2 \times 0.1, 0.8 \times 0.8 \times 0.1, 0.0 \times 0.2 \times 0.1) & = & \max(0.018 \times 0.1 \times 0.4, 0.0005 \times 0.2 \times 0.4, 0.0 \times 0.1 \times 0.4, 0.0 \times 0.0 \times 0.4) & = & \max(0.018 \times 0.1 \times 0.3, 0.0005 \times 0.2 \times 0.3, 0.0 \times 0.1 \times 0.3, 0.0 \times 0.1 \times 0.3) & = \\ \max(0.0018 \times 0.1 \times 0.4, 0.00005 \times 0.2 \times 0.4, 0.0 \times 0.1 \times 0.4, 0.0 \times 0.0 \times 0.4) & = & \max(0.0018 \times 0.1 \times 0.3, 0.00005 \times 0.2 \times 0.3, 0.0 \times 0.1 \times 0.3, 0.0 \times 0.1 \times 0.3) & = & \max(0.0018 \times 0.1 \times 0.9, 0.00005 \times 0.2 \times 0.9, 0.0 \times 0.1 \times 0.9, 0.0 \times 0.1 \times 0.9) & = \\ 0.0018 & = & 0.00005 & = & 0.000005 & = \end{aligned}$$

States and deltas over time:

$\delta_1(i)$	state
A	/
B	/
H	/
S	/

$\delta_2(i)$	state
A	/
B	/
H	/
S	/

$\delta_3(i)$	state
A	/
B	/
H	/
S	/

$\delta_4(i)$	state
A	/
B	/
H	/
S	/

Backtracking gives two answers:

$$X_{1:t} = \{A, A, A, A\} \text{ and}$$

$$X_{1:t} = \{A, A, S, A\}$$

- Naive Bayes assumes all features are independent effects of the label.
- Naive Bayes for text: $P(Y, W_1, \dots, W_n) = P(Y) \prod_i P(W_i | Y)$, where W_i is the word at position i , $Y \in \{\text{spam, ham}\}$.

Maximum likelihood estimation

- Given the observed set D , find θ to maximize the probability of D

$$\theta = \underset{\theta}{\operatorname{argmax}} P(D|\theta)$$

- Set the derivative of $P(D|\theta)$ with respect to θ to zero, and solve for θ .

- pretend that we have seen every outcome k extra times.

$$P_{\text{Lap},k} = \frac{c(x)+k}{N+k|X|}$$

Linear classifiers

feature values(inputs), weights (learned), activation (sum,

$$\text{activation}_w(x) = \sum_i w_i \phi_i(x))$$

Binary perceptron learning process

- start with weights=0

- for each training instance (x, y^*) : classify with current weights, no change if correct else adjust the weight vector by adding or subtracting the feature vector (subtract if $y^*=-1$).

Multiclass perceptron learning process

- start with weights=0
- pick up training examples one by one
- predict with current weights $\hat{y} = \text{argmax}_y(w_y \cdot \phi(x))$, no change if correct. Otherwise, lower score of wrong answer and raise score of right answer $w_{\hat{y}} = w_{\hat{y}} - \phi(x)$, $w_{y^*} = w_{y^*} + \phi(x)$

Probabilistic Perceptron

- softmax function

initialize w (e.g., randomly)

$$\frac{d}{dw_y} \log P_w(y_i|x_i) = x_i(I(y=y_i) - P(y|x_i))$$

repeat for K iterations:

for each example (x_i, y_i) :

compute gradient $\Delta_i = -\nabla_w \log P_w(y_i|x_i)$

compute gradient $\nabla_w \mathcal{L} = \sum_i \Delta_i$

$w \leftarrow w - \alpha \nabla_w \mathcal{L}$

❖ α : learning rate -- hyperparameter that needs to be chosen carefully

❖ How? Try multiple choices

❖ Crude rule of thumb: update should change w by about 0.1-1%

❖ False \models True ✓

❖ True \models False F

❖ $(A \wedge B) \models (A \Leftrightarrow B)$ ✓

❖ $(A \vee B) \wedge (\neg C \vee \neg D \vee E) \models (A \vee B)$ ✓

❖ $(A \vee B) \wedge \neg(A \Rightarrow B)$ is satisfiable ✓

❖ $(A \Leftrightarrow B) \wedge (\neg A \vee B)$ is satisfiable ✓

❖ $(A \Leftrightarrow B) \Leftrightarrow C$ has the same number of models as $(A \Leftrightarrow B)$ for any fixed set of proposition symbols that includes A, B, C ✓

3. (5 points) Use Bayes Rule to compute the posterior probabilities $\mathbb{P}(\text{spam} | \text{"viagra"})$ and $\mathbb{P}(\neg \text{spam} | \text{"viagra"})$.

Solution: The un-normalized posterior probabilities are:

$$\begin{aligned} \mathbb{P}(\text{spam} | \text{"viagra"}) &= \alpha \cdot \mathbb{P}(\text{"viagra"} | \text{spam}) \cdot \mathbb{P}(\text{spam}) = \alpha \cdot .015 \cdot .60 = \alpha \cdot .0090 \\ \mathbb{P}(\neg \text{spam} | \text{"viagra"}) &= \alpha \cdot \mathbb{P}(\text{"viagra"} | \neg \text{spam}) \cdot \mathbb{P}(\neg \text{spam}) = \alpha \cdot .001 \cdot .40 = \alpha \cdot .0004 \end{aligned}$$

Normalizing, so they add up to 1.0, we get: $\mathbb{P}(\text{spam} | \text{"viagra"}) = 90/94 = 0.9574$
 $\mathbb{P}(\neg \text{spam} | \text{"viagra"}) = 4/94 = 0.0426$.

4.2 Binary Classification

Recall a probabilistic binary classifier $f : \mathcal{X} \rightarrow \mathcal{Y}$ is defined as follows:

$$f(\mathbf{x}) = \begin{cases} 1 & \text{if } \eta \geq \frac{1}{2} \\ 0 & \text{if } \eta < \frac{1}{2} \end{cases}$$

where \mathbf{x} is known to be a non-negative vector and the two classes are encoded 0 and 1. We consider here a probabilistic binary classifier where η is defined by:

$$\eta_{\theta} = \frac{\mathbf{w}_1^T \mathbf{x}}{\mathbf{w}_0^T \mathbf{x} + \mathbf{w}_1^T \mathbf{x}}$$

where $\theta = (\mathbf{w}_0, \mathbf{w}_1) \in \mathbb{R}_{+}^{2(D+1)}$ and assuming each instance $\mathbf{x} \in \mathbb{R}^{D+1}$ has been appended a constant 1.

1. (3 points) Write the equation of the decision boundary for this classifier.

Solution: $(\mathbf{w}_0 - \mathbf{w}_1)^T \mathbf{x} = 0$

2. (3 points) What are the possible issues with this model?

Solution: The model defines a linear classifier. The weights may become negative while training, especially if the classification problem is not linearly separable.

3. We would like to learn parameters θ from a dataset $\mathcal{D} = \{(\mathbf{x}^i, y^i) | i = 1, \dots, N\}$, which has been collected in an i.i.d. way.

(a) (3 points) Write the log-likelihood of \mathcal{D} given parameter θ : $L(\theta; \mathcal{D}) = \log \mathbb{P}(\mathcal{D} | \theta)$.

$$\begin{aligned} \text{Solution: } \mathbb{P}(\mathcal{D} | \theta) &= \prod_{i=1}^N \left(\frac{\mathbf{w}_1^T \mathbf{x}^i}{\mathbf{w}_0^T \mathbf{x}^i + \mathbf{w}_1^T \mathbf{x}^i} \right)^{y^i} \left(\frac{\mathbf{w}_0^T \mathbf{x}^i}{\mathbf{w}_0^T \mathbf{x}^i + \mathbf{w}_1^T \mathbf{x}^i} \right)^{1-y^i} \\ L(\theta; \mathcal{D}) &= \log \mathbb{P}(\mathcal{D} | \theta) = \sum_{i=1}^N y^i \log \mathbf{w}_1^T \mathbf{x}^i + (1-y^i) \log \mathbf{w}_0^T \mathbf{x}^i - \log ((\mathbf{w}_0 + \mathbf{w}_1)^T \mathbf{x}^i) \end{aligned}$$

(b) (4 points) Write the gradient of this log-likelihood with respect to θ . Recall the gradient is given by:

$$\nabla_{\theta} L(\theta; \mathcal{D}) = \left(\left(\frac{\partial L(\theta; \mathcal{D})}{\partial w_{0k}} \right)_{k=1, \dots, D+1}, \left(\frac{\partial L(\theta; \mathcal{D})}{\partial w_{1k}} \right)_{k=1, \dots, D+1} \right)$$

$$\begin{aligned} \text{Solution: } \frac{\partial L(\theta; \mathcal{D})}{\partial w_{0k}} &= \sum_{i=1}^N \frac{y^i}{\mathbf{w}_0^T \mathbf{x}^i} x_k^i + \frac{1}{(\mathbf{w}_0 + \mathbf{w}_1)^T \mathbf{x}^i} x_k^i \\ \frac{\partial L(\theta; \mathcal{D})}{\partial w_{1k}} &= \sum_{i=1}^N \frac{1-y^i}{\mathbf{w}_1^T \mathbf{x}^i} x_k^i + \frac{1}{(\mathbf{w}_0 + \mathbf{w}_1)^T \mathbf{x}^i} x_k^i \end{aligned}$$

5 General Questions

For each question, answer "yes" or "no" and **justify** your answer.

1. (2 points) For any events E and F , we always have $\mathbb{P}(E | F) \geq \min(\mathbb{P}(E), \mathbb{P}(F))$

2. (2 points) For any events E and F , we always have $\mathbb{P}(E | F) \leq \max(\mathbb{P}(E), \mathbb{P}(F))$

Solution: (1pt for answer, 1 pt for justification)

No, for instance, for $E = F$ and $\mathbb{P}(E) < 1$, we have $\mathbb{P}(E | F) = 1$.

3. (2 points) We are given a joint distribution $\mathbb{P}(X, Y)$ over variables X and Y whose domains are respectively \mathcal{D}_X and \mathcal{D}_Y . We always have $\sum_{x \in \mathcal{D}_X} \mathbb{P}(X = x | Y = y) = 1$ for any $y \in \mathcal{D}_Y$ such that $\mathbb{P}(y) > 0$.

Solution: (1pt for answer, 1 pt for justification)

Yes, $\mathbb{P}(X | Y = y)$ is a probability distribution.

4. (2 points) We are given a joint distribution $\mathbb{P}(X, Y)$ over variables X and Y whose domains are respectively \mathcal{D}_X and \mathcal{D}_Y . Assume that $\forall y \in \mathcal{D}_Y, \mathbb{P}(y) > 0$. Then, we always have $\sum_{x \in \mathcal{D}_X} \mathbb{P}(X = x | Y = y) = 1$ for any $x \in \mathcal{D}_X$.

Solution: (1pt for answer, 1 pt for justification)

No, $\mathbb{P}(X = x | Y)$ is not a probability distribution.