

CSE 575: Statistical Machine Learning (Spring 2021)

Instructor: Nupur Thakur

Midterm Review

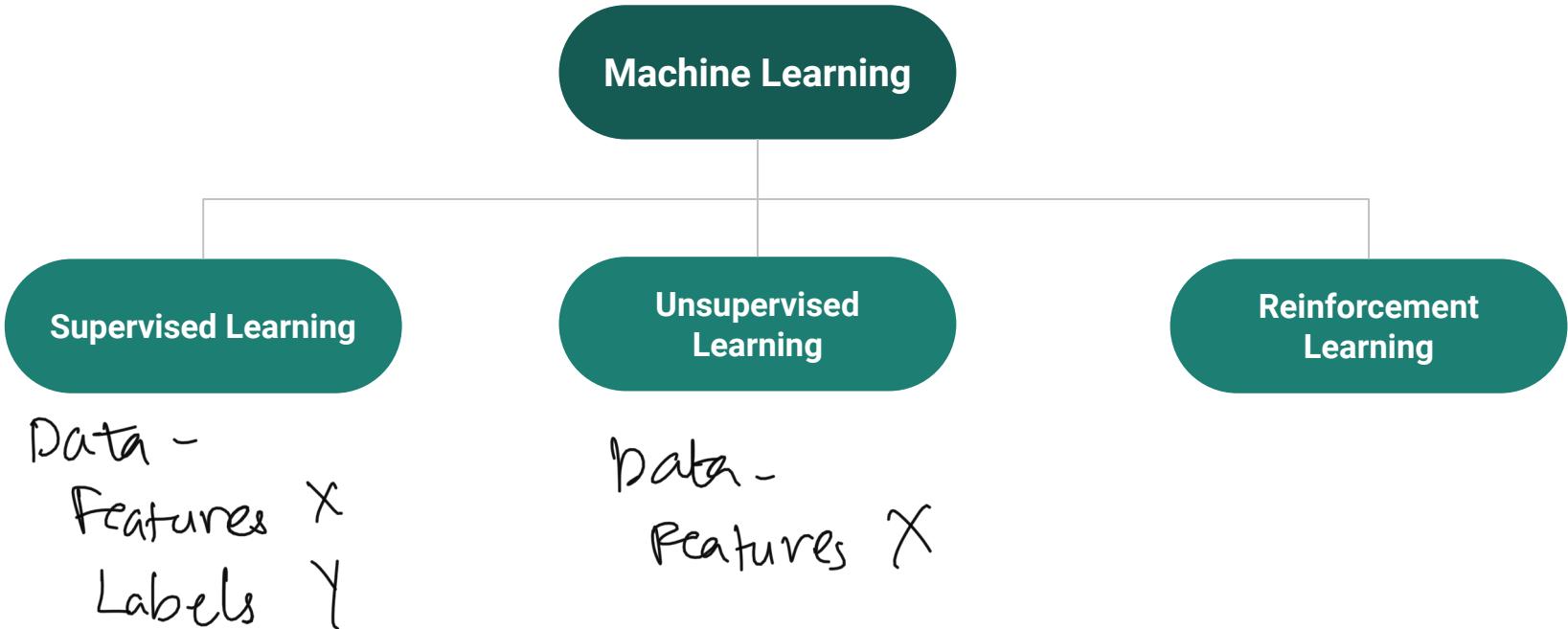


Table of contents

1. Concepts Review

2. Sample Questions

Types of Machine Learning



$$\text{Bias} - \underbrace{E[\hat{\theta}]}_{\text{Estimated Parameter}} - \theta = 0 \quad E[\hat{\theta}] = \theta$$

Maximum Likelihood Estimation

$$\underline{\underline{Ex}} - \underline{\underline{E[\hat{\mu}]}} = \mu \Rightarrow \text{Unbiased}$$

$$\underline{\underline{\hat{\mu}}} \neq \mu \Rightarrow \text{biased.}$$

$E[\mu] = \int x p(x) dx$ unbiased

- Given some training data and assuming a parametric model $\underline{\underline{p(x|\theta)}}$; what specific θ will fit/explain the data best?
- To consider all the samples denoted by $D = \{x_1, x_2, \dots, x_n\}$, assume that all the samples are i.i.d - independent and identically distributed.
- So, data likelihood represented by $L(\theta)$ is -

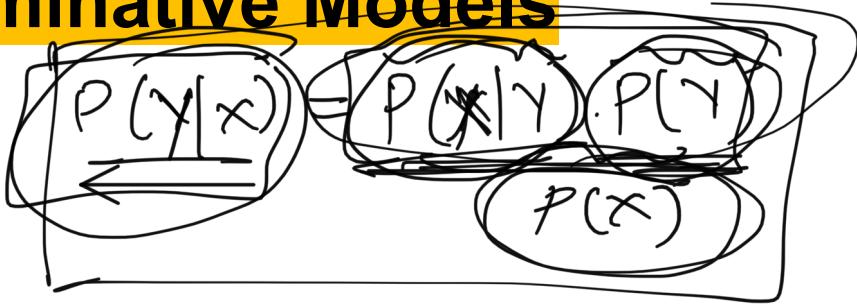
$$L(\theta) = P(D|\theta) = \prod_i P(x_i|\theta)$$

$$\boxed{\hat{\theta} = \operatorname{argmax}_{\theta} p(x|\theta)}$$

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$$

Generative vs Discriminative Models

- Generative models -
 - Learn $P(y)$ and $P(x|y)$.
 - Ex: Bayesian classifier, Naive Bayes.
- Discriminative models -
 - Directly learn $P(y|x)$
 - Ex: Logistic Regression



$$P(y|x)$$

Naive Bayes

Supervised

$$\mathbf{x} = \{x_1, \dots, x_n\}$$
$$y = 1$$

- The "naive" conditional independence assumption: each feature is (conditionally) independent of every other feature, given the label, i.e.,
 $p(x_i | \{x_j \text{ for any } j \neq i\}, y) = p(x_i | y)$
- The predicted label is given by -

$$\hat{y} = \operatorname{argmax}_y P(y) \prod_{i=1}^d p(x_i | y)$$

Predicted label.

Linear Regression

- Regression - A training set of n samples $\langle \underline{x}^{(i)}, \underline{y}^{(i)} \rangle$ where $\underline{y}^{(i)}$ is a continuous “label” (or target value) for $\underline{x}^{(i)}$
- Linear regression - modeling the relation between y and x via a linear function

$$y \approx w_0 + w_1 x_1 + \dots + w_d x_d = \mathbf{w}^t \mathbf{x}$$

- The error is given as - $\|e\|^2 = \|y - \underline{\mathbf{X}^t \mathbf{w}}\|^2$

Logistic Regression

$$P(y|x) = \text{logistic function} \\ w^T x \\ \equiv \\ w^T x \geq 0$$

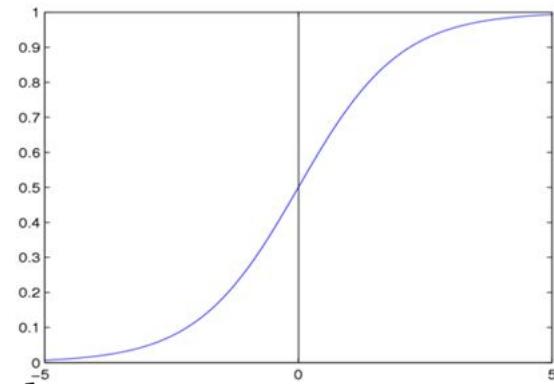
- Training set: n labelled samples $\langle \underline{x}(i), \underline{y}(i) \rangle$
- Use the logistic function for modeling $P(y|x)$, considering only the case of $y \in \{0,1\}$

$$P(y=0|x) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^d w_i x_i)}$$

$$P(y=1|x) = \frac{\exp(w_0 + \sum_{i=1}^d w_i x_i)}{1 + \exp(w_0 + \sum_{i=1}^d w_i x_i)}$$

→ Gradient ascent

$$w_{t+1} = w_t + \alpha \frac{\partial L(\omega)}{\partial \omega}$$



$$\sigma(t) = \frac{1}{1+e^{-t}} = \frac{e^t}{1+e^t}$$

Support Vector Machines (SVM)



- Key idea - To find the decision boundary such that the margin is maximized.
- Data - $\langle x^{(i)}, y^{(i)} \rangle$, $y^{(i)} \in \{-1, 1\}$, $x^{(i)} \in R^d$, for all $i=1, \dots, n$
- Plane equations-

$$\begin{cases} w^T x + b = 1 \\ w^T x + b = -1 \end{cases} \quad \text{Margin hyperplane.}$$

- Margin - $w^T x + b = 0$ — Decision boundary.

$$d = \frac{2}{\|w\|}$$

SVM - Problem Formulation

$$\{\mathbf{w}^*, b^*\} = \underset{\mathbf{w}, b}{\operatorname{argmin}} \|\mathbf{w}\| \text{ or } \{\mathbf{w}^*, b^*\} = \underset{\mathbf{w}, b}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{w}\|^2$$

Subject to

$$\begin{cases} \mathbf{w}^t \mathbf{x}^{(i)} + b \geq 1 & \text{for } y^{(i)} = +1 \\ \mathbf{w}^t \mathbf{x}^{(i)} + b \leq -1 & \text{for } y^{(i)} = -1 \end{cases}$$

The constraints can be combined into:

$$y^{(i)}(\mathbf{w}^t \mathbf{x}^{(i)} + b) - 1 \geq 0 \quad \forall i$$

Soft-Margin SVM Formulation

subject to

$$\{\mathbf{w}^*, b^*\} = \underset{\mathbf{w}, b}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{w}\|^2 + C(\sum \xi_i)$$

Parameter C controls the penalty.

slack variable

Very high values \rightarrow Hard margin SVM

$$\mathbf{w}^t \mathbf{x}^{(i)} + b \geq 1 - \xi_i \text{ for } y^{(i)} = +1$$

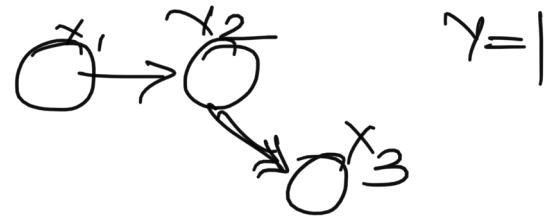
$$\mathbf{w}^t \mathbf{x}^{(i)} + b \leq -1 + \xi_i \text{ for } y^{(i)} = -1$$

$$\xi_i \geq 0, \forall i$$

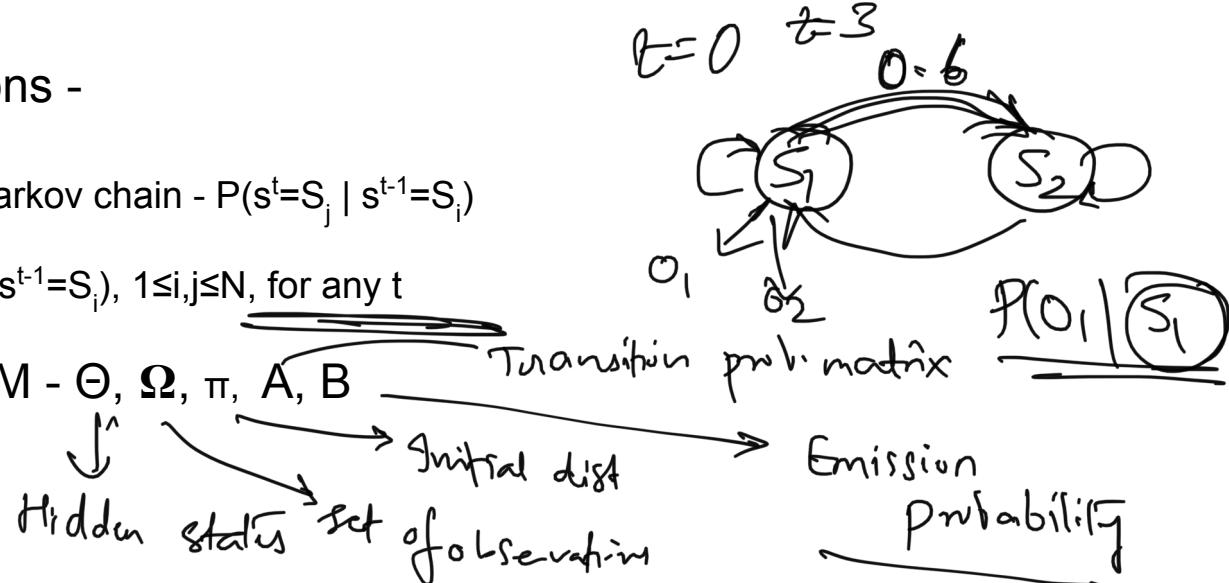
Very low values \rightarrow Larger misclassifications

$$\phi(x_i - x_j)$$

Hidden Markov Model (HMM)



- Dynamic Bayesian network - modeling the process indexed by time.
- Two assumptions -
 - First-order Markov chain - $P(s^t=S_j | s^{t-1}=S_i)$
 - $a_{ij} = P(s^t=S_j | s^{t-1}=S_i)$, $1 \leq i, j \leq N$, for any t
- Specifying HMM - $\Theta, \Omega, \pi, A, B$



Problems in HMM

- For a given HMM $\Lambda = \{\Theta, \Omega, A, B, \pi\}$
 - Estimation of model parameters
 - Given an observation sequence $O = \{o^1, o^2, \dots, o^k\}$, what is the most likely state sequence $S = \{s^1, s^2, \dots, s^k\}$ that has produced O ? $\xrightarrow{\text{Decoding}}$
 - How likely is an observation O ? $\xrightarrow{P(O)}$

Trained
HMM

HMM Parameter Estimation

- Given labeled data - state and observation

$$\text{A} \rightarrow t(S_i|S_j) = \frac{\text{number of } (s^t = S_i, s^{t-1} = S_j)}{\text{number of } S_j}$$

$$\text{B} \rightarrow e(o_r|S_j) = \frac{\text{number of } (o^t = o_r, s^t = S_j)}{\text{number of } S_j}$$

- Observation sequence — Forward-Backward Algo.

HMM State Estimation

- Given an observation (sequence) $O = \{o_1, o_2, \dots, o_k\}$, what is the most likely state sequence $S = \{s_1, s_2, \dots, s_k\}$ that has produced O ?

Viterbi

Initialization

$$\delta_{S_i}(1) = t(S_i|s^*)e(o^1|S_i), \quad \forall S_i \in \Theta$$

Induction:

for $2 \leq t \leq k$, do

$$\delta_{S_i}(t) = \max_{S_j} t(S_i|S_j)e(o^t|S_i)\delta_{S_j}(t-1)$$

$$\psi_{S_i}(t) = \operatorname{argmax}_{S_j} t(S_i|S_j)e(o^t|S_i)\delta_{S_j}(t-1)$$

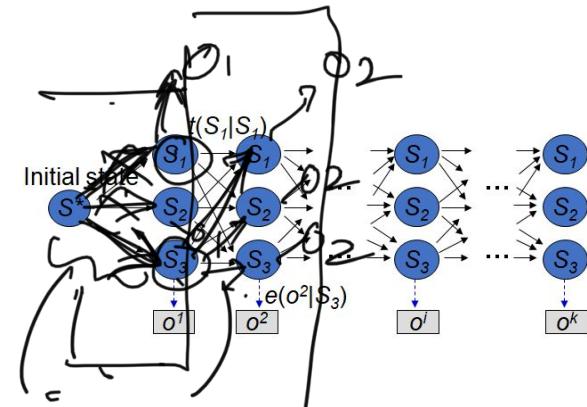
Termination:

- The probability of the best state sequence: $\max_{S_j} \delta_{S_j}(k)$

- The best last state: $\hat{s}^k = \operatorname{argmax}_{S_j} \delta_{S_j}(k)$

- Back trace to get other states:

$$\hat{s}^t = \psi_{\hat{s}^{t+1}}(t), \text{ for } t = k-1, \dots, 1.$$



s^*, s_3, s_1

HMM - Evaluate P(O)

- How likely is an observation O?

$$P(O) = \sum_s P(S, O)$$

s

- Forward algorithm-



Initialization: $\alpha_{S_i}(1) = t(S_i|s^*)e(o^1|S_i), \quad \forall S_i \in \Theta$

Induction:

for $2 \leq t \leq k$, do $\alpha_{S_i}(t) = \sum_{S_j} t(S_i|S_j)e(o^t|S_i)\alpha_{S_j}(t-1)$

Termination:

$$P(O) = \sum_{S_j} \alpha_{S_j}(k)$$

Sample Questions

Problem 1:

$$x_1 = 0 \quad y = 1$$

$$x_1 = 1 \quad y = 0$$

$$\{ P(x_1=1|y=1) \} \\ P(x_1=1|y=0)$$

$$1 - \boxed{P(x_1=0|y=1)} \quad \boxed{P(x_1=0|y=0)}$$

The following data is used for training Naive Bayes binary classifier. The last column is the binary class label; Each of the first 4 columns is a binary feature, and each row is a training example.

X1	X2	X3	X4	Y
1	0	0	0	1
0	1	1	0	1
1	0	0	1	0
0	1	1	1	0
1	1	1	0	1

$$P(x_2=0|y=1) = \frac{1}{3}$$
$$P(y=0) = \frac{2}{5}$$
$$P(x_1=1|y=1) = \frac{1}{3}$$
$$P(x_1=0|y=1) = \frac{2}{3}$$
$$P(x_1=0|y=0) = \frac{4}{5}$$

Problem 2:

State true/false along with justification:

1. Logistic Regression may give us a non-linear classifier, depending on how the training examples are distributed in the feature space.

Ans - False because logistic regression is a linear classifier, producing a linear decision boundary.

Problem 3:

Multiple Choice Questions:

Which of the following statements is true about HMM?

- a. Given some initial state and transition matrix, it is possible that there exists a state that you can never achieve.
- b. The elements of a column of the state transition matrix always sum to 1.
- c. The sum of a row of the transition matrix may not be 1.

Ans ~ (9)