



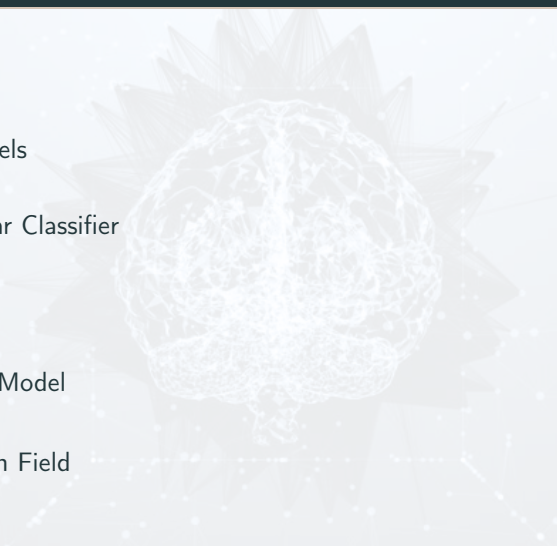
Machine Learning

Discriminative Models

Lecturer: Duc Dung Nguyen, PhD.

Contact: nddung@hcmut.edu.vn

Faculty of Computer Science and Engineering
Hochiminh city University of Technology

- 
- A large, faint, stylized graphic of a brain with neural connections, overlaid with a starburst pattern, serving as a background for the list.
1. Discriminative Models
 2. Feature-based Linear Classifier
 3. Logistic Regression
 4. Maximum Entropy Model
 5. Conditional Random Field

Discriminative Models

Generative model (BNs, HMMs)

- Joint distributions: $p(\mathbf{y}, \mathbf{x})$
- It can generate any distribution on \mathbf{y} and \mathbf{x} .

Discriminative model:

- Conditional distributions: $p(\mathbf{y}|\mathbf{x})$
- It discriminates \mathbf{y} given \mathbf{x} .

Feature-based Linear Classifier

- Decision is based on a linear combination of features.

- Decision is based on a linear combination of features.
- Features are essential pieces of information in an observation that decide classification.

- Decision is based on a linear combination of features.
- Features are essential pieces of information in an observation that decide classification.
- A feature can be defined as a function with a bounded real value of a class and an observation:

$$f(c, \mathbf{x})$$

- Decision is based on a linear combination of features.
- Features are essential pieces of information in an observation that decide classification.
- A feature can be defined as a function with a bounded real value of a class and an observation:

$$f(c, \mathbf{x})$$

- Linear classifier:

$$\arg \max_{c \in C} \sum_{m=1..M} \lambda_m \cdot f_m(c, x)$$

- SVMs:

$$y(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + \mathbf{b} = \sum_{m=1..M} \mathbf{w}_m \cdot \mathbf{x}_m + \mathbf{b}$$

Feature function: $f_m(c, \mathbf{x}) = x_m$

- Naive Bayes classifier:

$$c_{NB} = \arg \max_{c \in C} p(c) \cdot \prod_{m=1}^M p(x_m|c)$$

$$c_{NB} = \arg \max_{c \in C} \left(\log p(c) + \sum_{m=1}^M \log p(x_m|c) \right)$$

Feature function:

$$f_m(c, \mathbf{x}) = \log p(x_m|c)$$

- In NLP, a feature could be an indicator function whose value is boolean:

$$f(c, \mathbf{x}) = (c \text{ is a certain class}) \wedge (\mathbf{x} \text{ has a certain property})$$

- In NLP, a feature could be an indicator function whose value is boolean:

$$f(c, \mathbf{x}) = (c \text{ is a certain class}) \wedge (\mathbf{x} \text{ has a certain property})$$

- Example:

$$f(c, \mathbf{x}) = (c = \text{NUMBER}) \text{ and } (\mathbf{x} \text{ contains only digits})$$

$$f(\text{NUMBER}, "2018") = 1$$

$$f(\text{NUMBER}, "may 2018") = 0$$

$$f(\text{IDENTIFIER}, "may 2018") = 0$$

- Empirical count of a feature:

$$\tilde{E}(f_m) = \sum_{(c, \mathbf{x}) \in \text{observed}(c, \mathbf{x})} \tilde{p}(c, \mathbf{x}) \cdot f_m(c, \mathbf{x})$$

$$\tilde{E}(f_m) = \frac{1}{D} \sum_{(c, \mathbf{x}) \in \text{observed}(c, \mathbf{x})} f_m(c, \mathbf{x})$$

- Model expectation of a feature:

$$E(f_m) = \sum_{(c, \mathbf{x}) \in (C, X)} p(c, \mathbf{x}) \cdot f_m(c, \mathbf{x})$$

$$E(f_m) = \sum_{(c, \mathbf{x}) \in (C, X)} p(\mathbf{x}) \cdot p(c|\mathbf{x}) \cdot f_m(c, \mathbf{x})$$

$$E(f_m) \approx \sum_{(c, \mathbf{x}) \in (C, X)} \tilde{p}(\mathbf{x}) \cdot p(c|\mathbf{x}) \cdot f_m(c, \mathbf{x})$$

$$E(f_m) \approx \frac{1}{D} \sum_{\mathbf{x} \in \text{observed}(\mathbf{x})} \sum_{c \in C} p(c|\mathbf{x}) \cdot f_m(c, \mathbf{x})$$

- Empirical count of a feature:

$$\tilde{E}(f_m) = \frac{1}{D} \sum_{(c, \mathbf{x}) \in \text{observed}(c, \mathbf{x})} f_m(c, \mathbf{x})$$

- Model expectation of a feature:

$$E(f_m) = \frac{1}{D} \sum_{\mathbf{x} \in \text{observed}(\mathbf{x})} \sum_{c \in C} p(c|\mathbf{x}) \cdot f_m(c, \mathbf{x})$$

- Consistency constraint:

$$E(f_m) = \tilde{E}(f_m)$$

Logistic Regression

- A discriminative model:

$$p(y|\mathbf{x}) = \frac{\exp \sum_{m=1..M} \lambda_m \cdot f_m(y, \mathbf{x})}{\sum_{y' \in Y} \exp \sum_{m=1..M} \lambda_m \cdot f_m(y', \mathbf{x})}$$

- A discriminative model:

$$p(y|\mathbf{x}) = \frac{\exp \sum_{m=1..M} \lambda_m \cdot f_m(y, \mathbf{x})}{\sum_{y' \in Y} \exp \sum_{m=1..M} \lambda_m \cdot f_m(y', \mathbf{x})}$$

- Not factorized into a product of conditional distributions, but a product of arbitrary functions.

- A discriminative model:

$$p(y|\mathbf{x}) = \frac{\exp \sum_{m=1..M} \lambda_m \cdot f_m(y, \mathbf{x})}{\sum_{y' \in Y} \exp \sum_{m=1..M} \lambda_m \cdot f_m(y', \mathbf{x})}$$

- Not factorized into a product of conditional distributions, but a product of arbitrary functions.
- Linear classifier:

$$\log p(y|\mathbf{x}) = \sum_{m=1..M} \lambda_m \cdot f_m(y, \mathbf{x})$$

- Comparison to Naive Bayes classifier:

$$\begin{aligned}y_{NB} &= \arg \max_{y \in Y} \prod_{m=1..M} p(x_m|y) \cdot p(y) \\&= \arg \max_{y \in Y} \exp(\log p(y) + \sum_{m=1..M} \log p(x_m|y)) \\&= \arg \max_{y \in Y} \exp \sum_{m=1..M} \lambda_m f_m(y, x_m)\end{aligned}$$

- Comparison to Naive Bayes classifier:

$$\begin{aligned}y_{NB} &= \arg \max_{y \in Y} \prod_{m=1..M} p(x_m|y) \cdot p(y) \\&= \arg \max_{y \in Y} \exp(\log p(y) + \sum_{m=1..M} \log p(x_m|y)) \\&= \arg \max_{y \in Y} \exp \sum_{m=1..M} \lambda_m f_m(y, x_m)\end{aligned}$$

- Naive Bayes classifier is just an exponential model.

- A discriminative model:

$$p(y|\mathbf{x}) = \frac{\exp \sum_{m=1..M} \lambda_m \cdot f_m(y, \mathbf{x})}{\sum_{y^* \in Y} \exp \sum_{m=1..M} \lambda_m \cdot f_m(y^*, \mathbf{x})}$$

- A two-class case:

$$p(\oplus|\mathbf{x}) = \frac{\exp \sum_{m=1..M} \lambda_m \cdot f_m(y, \mathbf{x})}{1 + \exp \sum_{m=1..M} \lambda_m \cdot f_m(y, \mathbf{x})} p(\ominus|\mathbf{x}) = \frac{1}{1 + \exp \sum_{m=1..M} \lambda_m \cdot f_m(y, \mathbf{x})}$$

- Learning the parameters λ 's by an iterative algorithm until the convergence of the consistency constraint.

Maximum Entropy Model

- **Principle of maximum entropy:** the only unbiased assumption is a distribution that is as uniform as possible given the available information.

- **Principle of maximum entropy:** the only unbiased assumption is a distribution that is as uniform as possible given the available information.
- **Maximum entropy model:** the proper probability distribution is the one that maximizes the entropy given the constraints from the training data.

- Conditional entropy:

$$H(y|\mathbf{x}) = - \sum_{(y,\mathbf{x}) \in (Y,X)} p(y,\mathbf{x}) \cdot \log p(y|\mathbf{x})$$

$$H(y|\mathbf{x}) \approx - \sum_{(y,\mathbf{x}) \in (Y,X)} \tilde{p}(\mathbf{x}) \cdot p(y|\mathbf{x}) \cdot \log p(y|\mathbf{x})$$

- Maximum entropy model:

$$p^*(y|\mathbf{x}) = \arg \max_{p(y|\mathbf{x}) \in P} H(y|\mathbf{x})$$

Constraints:

$$E(f_m) = \tilde{E}(f_m)$$

$$\sum_{y \in Y} p(y|\mathbf{x}) = 1$$

- Lagrange function $L(p(y|\mathbf{x}))$:

$$H(y|\mathbf{x}) + \sum_{m=1..M} \lambda_m (E(f_m) - \tilde{E}(f_m)) + \lambda_{M+1} (\sum_{y \in Y} p(y|\mathbf{x}) - 1)$$

- Optimization:

$$\partial L(p(y|\mathbf{x})) / \partial p(y|\mathbf{x}) = 0$$

$$H(y|\mathbf{x}) = - \sum_{(y, \mathbf{x}) \in (Y, X)} \tilde{p}(\mathbf{x}) \cdot p(y|\mathbf{x}) \cdot \log p(y|\mathbf{x})$$

$$\partial H(y|\mathbf{x}) / \partial p(y|\mathbf{x}) = -\tilde{p}(\mathbf{x}) \cdot (\log p(y|\mathbf{x}) + 1)$$

- Lagrange function $L(p(y|x))$:

$$H(y|\mathbf{x}) + \sum_{m=1..M} \lambda_m (E(f_m) - \tilde{E}(f_m)) + \lambda_{M+1} (\sum_{y \in Y} p(y|\mathbf{x}) - 1)$$

- Optimization:

$$E(f_m) - \tilde{E}(f_m) = \sum_{(y,x) \in (Y,X)} \tilde{p}(\mathbf{x}) \cdot p(y|\mathbf{x}) \cdot f_m(y, \mathbf{x}) - \sum_{(y,\mathbf{x}) \in (Y,X)} \tilde{p}(y, \mathbf{x}) \cdot f_m(y, \mathbf{x})$$

$$\partial \sum_{m=1..M} \lambda_m (E(f_m) - \tilde{E}(f_m)) / \partial p(y|\mathbf{x}) = \sum_{m=1..M} \lambda_m \tilde{p}(\mathbf{x}) \cdot f_m(y, \mathbf{x})$$

- Lagrange function $L(p(y|\mathbf{x}))$:

$$H(y|\mathbf{x}) + \sum_{m=1..M} \lambda_m (E(f_m) - \tilde{E}(f_m)) + \lambda_{M+1} (\sum_{y \in Y} p(y|\mathbf{x}) - 1)$$

- Optimization:

$$\partial \lambda_{M+1} (\sum_{y \in Y} p(y|\mathbf{x}) - 1) / \partial p(y|\mathbf{x}) = \lambda_{M+1}$$

- Lagrange function $L(p(y|\mathbf{x}))$:

$$H(y|x) + \sum_{m=1..M} \lambda_m (E(f_m) - \tilde{E}(f_m)) + \lambda_{M+1} (\sum_{y \in Y} p(y|\mathbf{x}) - 1)$$

- Optimization:

$$-\tilde{p}(\mathbf{x}) \cdot (\log p(y|\mathbf{x}) + 1) + \sum_{m=1..M} \lambda_m \tilde{p}(\mathbf{x}) \cdot f_m(y, \mathbf{x}) + \lambda_{M+1} = 0$$

$$p(y|\mathbf{x}) = \exp \sum_{m=1..M} \lambda_m f_m(y, \mathbf{x}) \exp(\lambda_{M+1} / \tilde{p}(\mathbf{x}) - 1)$$

$$\sum_{y \in Y} p(y|\mathbf{x}) = 1$$

- Lagrange function $L(p(y|\mathbf{x}))$:

$$H(y|\mathbf{x}) + \sum_{m=1..M} \lambda_m (E(f_m) - \tilde{E}(f_m)) + \lambda_{M+1} (\sum_{y \in Y} p(y|\mathbf{x}) - 1)$$

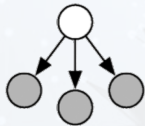
- Optimization:

$$\frac{\partial L(p(y|\mathbf{x}))}{\partial p(y|\mathbf{x})} = 0$$

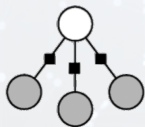
- Solution:

$$p(y|\mathbf{x}) = \frac{\exp \sum_{m=1..M} \lambda_m \cdot f_m(y, \mathbf{x})}{\sum_{y^* \in Y} \exp \sum_{m=1..M} \lambda_m \cdot f_m(y^*, \mathbf{x})}$$

Conditional Random Field



Naive Bayes



Logistic Regression



SEQUENCE



HMMs



Linear-chain CRFs

- Linear chain CRF:

$$p(\mathbf{y}|\mathbf{x}) = \frac{\prod_{t=1..T} (\exp \sum_{m=1..M} \lambda_m \cdot f_m(y_t, y_{t-1}, \mathbf{x}_t))}{\sum_{y^* \in Y} \prod_{t=1..T} (\exp \sum_{m=1..M} \lambda_m \cdot f_m(y_t^*, y_{t-1}^*, \mathbf{x}_t))}$$