# CSE 250A. Principles of AI

## Probabilistic Reasoning and Decision-Making

**Lecture 8 – Learning from Complete Data**

Lawrence Saul
Department of Computer Science and Engineering
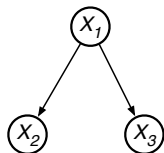University of California, San Diego

Fall 2021

# Outline

1 **Maximum likelihood**

2 **Markov models**

3 **Naive Bayes models**

4 **Preview**

## Learning in BNs (review)

**ASSUMPTIONS**

- Discrete random variables $\{X_1, X_2, \ldots, X_n\}$

- DAG is specified, assumed to be known and fixed.

- CPTs enumerate $P(X_i = x | \mathrm{pa}_i = \pi)$.

- IID data $\left\{ (x_1^{(t)}, x_2^{(t)}, \ldots, x_n^{(t)}) \right\}_{t=1}^{T}$



| example | $x_1$ | $x_2$ | $x_3$ |
|---------|-------|-------|-------|
| **1** | 1 | 4 | 5 |
| **2** | 3 | 2 | 4 |
| **3** | 2 | 1 | 3 |
| ⋮ | ⋮ | ⋮ | ⋮ |
| **T** | 1 | 3 | 2 |

Each example gives a **complete** instantiation of the nodes in the belief network.

## Computing the log-likelihood

$$
\begin{aligned}
\mathcal{L} &= \log P(\text{data}) \\
&= \log \prod_{t=1}^{T} P\left(x_1^{(t)}, x_2^{(t)}, \ldots, x_n^{(t)}\right) \qquad \boxed{\textbf{IID data}} \\
&= \log \prod_{t=1}^{T} \prod_{i=1}^{n} P\left(x_i^{(t)} \,\middle|\, \mathrm{pa}_i^{(t)}\right) \qquad \boxed{\textbf{product rule in BN}} \\
&= \sum_{i=1}^{n} \sum_{t=1}^{T} \log P\left(x_i^{(t)} \,\middle|\, \mathrm{pa}_i^{(t)}\right) \qquad \boxed{\textcolor{red}{\textbf{unweighted} \textbf{ sum over examples}}} \\
&= \sum_{i=1}^{n} \sum_{x} \sum_{\pi} \mathrm{count}(X_i = x, \mathrm{pa}_i = \pi) \, \log P(X_i = x | \mathrm{pa}_i = \pi)
\end{aligned}
$$

$$\boxed{\textcolor{blue}{\textbf{weighted} \textbf{ sum over co-occurrences}}}$$

## Interpreting the log-likelihood

$$\mathcal{L} = \sum_i \sum_x \sum_\pi \overbrace{\text{count}(X_i\!=\!x, \text{pa}_i\!=\!\pi)}^{\textbf{constants of the data}} \underbrace{\log P(X_i\!=\!x|\text{pa}_i\!=\!\pi)}_{\textbf{CPTs to optimize}}$$

- **The log-likelihood for complete data is a triple sum over**

  $i$ — the nodes in the BN

  $x$ — the values of each node $X_i$

  $\pi$ — the values $\pi$ of the parents of $X_i$

- **How to optimize?**

  Intuitively, the larger the $\text{count}(X_i\!=\!x, \text{pa}_i\!=\!\pi)$,
  the larger we should choose $P(X_i\!=\!x|\text{pa}_i\!=\!\pi)$.

# Decomposing the log-likelihood

- **Log-likelihood for BN**

$$\mathcal{L} \;=\; \sum_i \sum_\pi \sum_x \mathrm{count}(X_i\!=\!x, \mathrm{pa}_i\!=\!\pi) \, \log P(X_i\!=\!x|\mathrm{pa}_i\!=\!\pi)$$

- **Contribution from row $\pi$ of $i^{th}$ node's CPT**

$$\mathcal{L}_{i\pi} \;=\; \sum_x \mathrm{count}(X_i\!=\!x, \mathrm{pa}_i\!=\!\pi) \, \log P(X_i\!=\!x|\mathrm{pa}_i\!=\!\pi)$$

- **Divide and conquer**

  The overall optimization over $\mathcal{L}$ reduces to many simpler and smaller optimizations over each $\mathcal{L}_{i\pi}$.

---

*This is a special property of ML estimation for* **complete** *data.*

# ML Estimation

- **Problem**

  For each node $X_i$ in the BN, and for each row $\pi$ of its CPT, our goal is to maximize

  $$\mathcal{L}_{i\pi} = \sum_x \text{count}(X_i = x, \text{pa}_i = \pi) \log P(X_i = x | \text{pa}_i = \pi)$$

  subject to two constraints:

  1. $\sum_x P(X_i = x | \text{pa}_i = \pi) = 1$     *(normalized)*
  2. $P(X_i = x | \text{pa}_i = \pi) \geq 0$     *(nonnegative)*

- **Shorthand**

  $$\begin{aligned} C_\alpha &= \text{count}(X_i = \alpha, \text{pa}_i = \pi) \\ p_\alpha &= P(X_i = \alpha | \text{pa}_i = \pi) \end{aligned} \implies$$

  **How to maximize $\sum_\alpha C_\alpha \log p_\alpha$ such that $\sum_\alpha p_\alpha = 1$ and $p_\alpha \geq 0$?**

## Maximizing the likelihood

- **Compute the normalized counts:**

  Define $q_\alpha = \frac{C_\alpha}{\sum_\beta C_\beta}$ so that $\sum_\alpha q_\alpha = 1$ .

  Note that $q_\alpha$ is itself a distribution.

- **All these problems have the same solution:**

  | | | | |
  |---|---|---|---|
  | Maximize | $\sum_\alpha C_\alpha \log p_\alpha$ | such that | $\sum_\alpha p_\alpha = 1,\ p_\alpha \geq 0.$ |
  | **Minimize** | $\sum_\alpha C_\alpha \log \frac{1}{p_\alpha}$ | such that | $\sum_\alpha p_\alpha = 1,\ p_\alpha \geq 0.$ |
  | Minimize | $\sum_\alpha C_\alpha \log \frac{C_\alpha}{p_\alpha}$ | such that | $\sum_\alpha p_\alpha = 1,\ p_\alpha \geq 0.$ |
  | Minimize | $\underbrace{\sum_\alpha q_\alpha \log \frac{q_\alpha}{p_\alpha}}_{\mathrm{KL}(q,p)\ \leftarrow}$ | such that | $\sum_\alpha p_\alpha = 1,\ p_\alpha \geq 0.$ |

  KL distance from HW 1

  **Solution:** $p_\alpha = q_\alpha$

## ML solution from normalized counts

$$P_{\mathrm{ML}}(X_i\!=\!x|\mathrm{pa}_i\!=\!\pi) \;=\; \frac{\mathrm{count}(X_i\!=\!x, \mathrm{pa}_i\!=\!\pi)}{\sum_{x'} \mathrm{count}(X_i\!=\!x', \mathrm{pa}_i\!=\!\pi)}$$

- **For nodes with parents:**

$$P_{\mathrm{ML}}(X_i\!=\!x|\mathrm{pa}_i\!=\!\pi) \;=\; \frac{\mathrm{count}(X_i\!=\!x, \mathrm{pa}_i\!=\!\pi)}{\mathrm{count}(\mathrm{pa}_i\!=\!\pi)}$$

- **For root nodes:**

$$P_{\mathrm{ML}}(X_i\!=\!x) \;=\; \frac{\mathrm{count}(X_i\!=\!x)}{T}$$

## Properties of ML solution

- **Asymptotically correct:**

  The more data you have, the better your estimates.
  If $P(x_1, x_2, \ldots, x_n) > 0$, then

  $$\lim_{T \to \infty} P_{\mathrm{ML}}(x_1, x_2, \ldots, x_n) \;=\; P(x_1, x_2, \ldots, x_n)$$

- **But problematic for sparse data:**

  $$P_{\mathrm{ML}}(X_i \!=\! x | \mathrm{pa}_i \!=\! \pi) \;=\; \frac{\mathrm{count}(X_i \!=\! x, \mathrm{pa}_i \!=\! \pi)}{\mathrm{count}(\mathrm{pa}_i \!=\! \pi)}$$

  This is **undefined** when $\mathrm{count}(\mathrm{pa}_i \!=\! \pi) = 0$.
  Otherwise it is **zero** when $\mathrm{count}(X_i \!=\! x, \mathrm{pa}_i \!=\! \pi) = 0$.

# Outline

1. **Maximum likelihood estimation**

2. **Markov models**

3. **Naive Bayes models**

4. **Preview**

## Statistical language modeling

Let $w_\ell$ denote the $\ell^{\mathrm{th}}$ word in a sentence (or text).
How to model $P(w_1, w_2, \ldots, w_L)$?



**automatic speech recognition**



**machine translation**

## Context and expectations in language
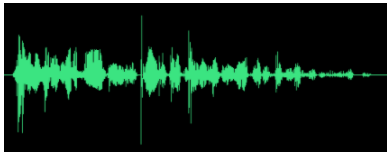


**"It's hard to wreck a nice beach."**



**"It's hard to recognize speech."**

# Simplifying assumptions

**❶ Finite context**

To predict the $\ell^{\text{th}}$ word, it is sufficient to consider a *finite* number of words that precede it:

$$P(w_\ell | w_1, w_2, \ldots, w_{\ell-1}) = P(w_\ell | \underbrace{w_{\ell-(n-1)}, \ldots, w_{\ell-1}}_{n-1 \text{ previous words}})$$

**❷ Position invariance**

Predictions should not depend on where the context occurs in the sentence or text:

$$P(W_\ell = w' | w_{\ell-(n-1)}, \ldots, w_{\ell-1})$$

$$= P(W_{s+\ell} = w' | W_{s+\ell-(n-1)} = w_{\ell-(n-1)}, \ldots, W_{s+\ell-1} = w_{\ell-1})$$

## Markov models

$P(w_1, w_2, \ldots, w_L)$

$$= \prod_\ell P(w_\ell | w_1, w_2, \ldots, w_{\ell-1}) \quad \boxed{\textbf{product rule}}$$

$$= \prod_\ell P(w_\ell | w_{\ell-(n-1)}, \ldots, w_{\ell-1}) \quad \boxed{\textbf{conditional independence}}$$

$\boxed{\textbf{Models of different orders}}$

$n = 1$ **unigram** $\quad (w_1) \quad (w_2) \quad (w_3) \quad \cdots \quad (w_{L-1}) \quad (w_L)$

$n = 2$ **bigram** $\quad (w_1) \rightarrow (w_2) \rightarrow (w_3) \quad \cdots \quad (w_{L-1}) \rightarrow (w_L)$

$n = 3$ **trigram**

## Markov models

$P(w_1, w_2, \ldots, w_L)$

$$= \prod_\ell P(w_\ell | w_1, w_2, \ldots, w_{\ell-1}) \quad \boxed{\textbf{product rule}}$$

$$= \prod_\ell P(w_\ell | w_{\ell-(n-1)}, \ldots, w_{\ell-1}) \quad \boxed{\textbf{conditional independence}}$$

$\boxed{\textbf{Models of different orders}}$



$n = 1$ **unigram** $\quad (w_1) \quad (w_2) \quad (w_3) \cdots (w_{L-1}) \quad (w_L)$

$n = 2$ **bigram** $\quad (w_1) \rightarrow (w_2) \rightarrow (w_3) \cdots (w_{L-1}) \rightarrow (w_L)$

$n = 3$ **trigram** $\quad (w_1) \rightarrow (w_2) \rightarrow (w_3) \cdots (w_{L-1}) \rightarrow (w_L)$

## Bigram models



Note that the same CPT
for $P(w_\ell = w' | w_{\ell-1} = w)$ is
used at each node (for $\ell > 1$).

**How to learn?**

**Collect** a large corpus of text with a well-defined vocabulary.

**Count** how often word $w$ is followed by the word $w'$.
**Count** how often word $w$ is followed by any word.

**Estimate** from empirical frequencies:

$$P_{\mathrm{ML}}(w_\ell = w' | w_{\ell-1} = w) \; = \; \frac{\mathrm{count}(w \rightarrow w')}{\mathrm{count}(w \rightarrow *)} \; = \; \frac{\mathrm{count}(w \rightarrow w')}{\sum_{w''} \mathrm{count}(w \rightarrow w'')}$$

## Problems with ML estimates

1. **No generalization to unseen $n$-grams:**

   ML estimates assign **zero** probability to $n$-grams that do not appear in the training corpus.
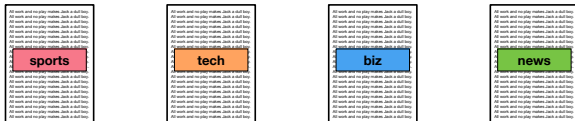
2. **The larger $n$, the worse the problem:**

   $n$-gram counts become increasingly sparse as $n$ increases. Many possible (but improbable) $n$-grams are not observed.

   **You will explore this problem further in HW 4.**

## Outline

1. **Maximum likelihood estimation**

2. **Markov models**

3. **Naive Bayes models**

4. **Preview**

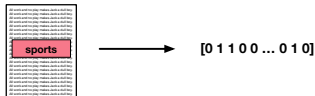# Document classification



- **Setup**

  Each document can be labeled by one of $m$ topics.
  Each document consists of words from a finite vocabulary.
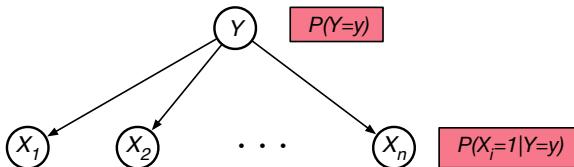
- **Random variables**

  Let $Y \in \{1, 2, \ldots, m\}$ denote the label.
  Let $X_i \in \{0, 1\}$ denote whether the $i^{\text{th}}$ word appears.

  This representation maps each document to a sparse binary vector of fixed length.



[ 0 1 1 0 0 ... 0 1 0 ]
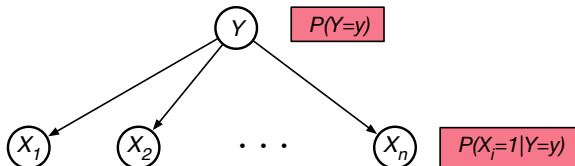
## Belief network



This DAG makes a fairly drastic assumption of conditional independence:

$$P(X_1, \ldots, X_n | Y) = \prod_{i=1}^{n} P(X_i | Y)$$

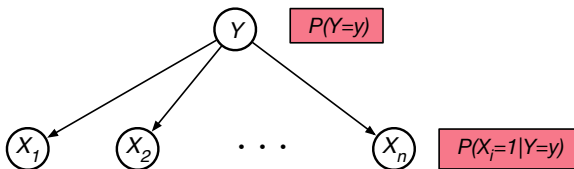For this reason it is called a **Naive Bayes** model.

## Naive Bayes model



Suppose this DAG is given, but the CPTs are not specified.
**How to learn the CPTs from data?**

- **Collect** a large corpus of documents.
- **Label** each document by a topic.
- **Estimate** the CPTs by maximizing the likelihood.
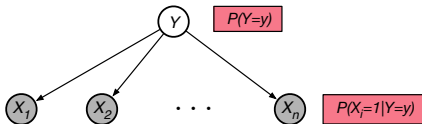
## ML estimation



$$P_{\mathrm{ML}}(Y=y) \quad = \quad \text{fraction of documents with label } y \text{ in the corpus}$$

$$P_{\mathrm{ML}}(X_i=1|Y=y) \quad = \quad \text{fraction of documents with label } y \text{ that contain the } i^{\mathrm{th}} \text{ word in the vocabulary}$$

**Once the model is learned, what is it good for?**

## Inference

**How to classify
an unlabeled
document?**



$P(Y=y|X_1, X_2, \ldots, X_n)$

$$= \frac{P(X_1, X_2, \ldots, X_n|Y=y)\,P(Y=y)}{P(X_1, X_2, \ldots, X_n)} \quad \boxed{\textbf{Bayes rule}}$$
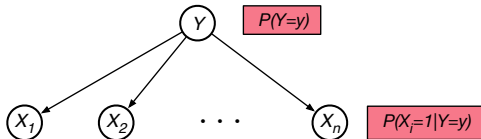
$$= \frac{P(Y=y)\prod_{i=1}^{n} P(X_i|Y=y)}{P(X_1, X_2, \ldots, X_n)} \quad \boxed{\textbf{conditional independence}}$$

$$= \frac{P(Y=y)\prod_{i=1}^{n} P(X_i|Y=y)}{\sum_{y'} P(Y=y')\prod_{i=1}^{n} P(X_i|Y=y')} \quad \boxed{\textbf{normalization}}$$

## Strengths and weaknesses



**Strengths**

- Easy to learn from data.
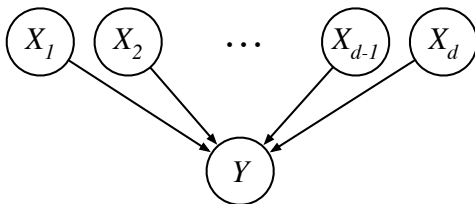- Easy to classify unlabeled documents.

**Weaknesses**

- Naive Bayes assumption of conditional independence
- No information about word ordering
- Binarization of word counts
- Etc ...

## Outline

**1** **Maximum likelihood estimation**

**2** **Markov models**

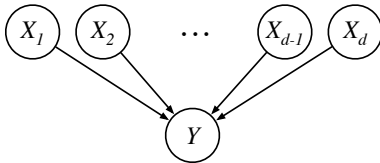**3** **Naive Bayes models**

**4** **Preview**

## Parametric models

**If the parent nodes are real-valued, then it is no longer possible to enumerate a conditional probability table.**



**How to predict $Y$ from real-valued parents $\vec{X} \in \mathbb{R}^d$?**

## Gaussian model



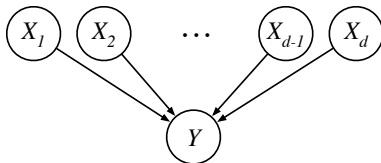Suppose $Y \in \mathbb{R}$ is a real-valued random variable.
Then we can use a **Gaussian conditional distribution**:

$$P\left(y|\vec{x}\right) \;=\; \frac{1}{\sqrt{2\pi\sigma^2}} \, \exp\left\{-\frac{(y - \vec{w} \cdot \vec{x})^2}{2\sigma^2}\right\}$$

How to learn the parameters $\sigma^2$ and $\vec{w} = (w_1, w_2, \ldots, w_d)$?
This is the problem of **linear regression**.

## Sigmoid model



Suppose $Y \in \{0, 1\}$ is a binary random variable.
Then we can use a **sigmoid conditional distribution**:

$$P(Y=1|\vec{x}) = \sigma(\vec{w} \cdot \vec{x}) = \frac{1}{1 + e^{-\vec{w} \cdot \vec{x}}}$$

How to learn the parameter $\vec{w} \in \mathbb{R}^d$?
This is the problem of **logistic regression**.

## Before next lecture ...

**Now would be a good time to review your linear algebra:**

- dot products
- matrix-vector multiplication
- systems of linear equations

**And also your multivariable calculus:**

- functions of several real variables
- partial derivatives
- gradients and Hessians