

# DSC250: Advanced Data Mining

## Machine Learning Basics

**Zhiting Hu**

Lecture 5, Jan 21, 2025

**UC San Diego**

**HALICIOĞLU DATA SCIENCE INSTITUTE**

# Outline

- Probability
  - Bayes' rule
  - Exponential family
  - Probabilistic graphical models
  - Entropy, KL divergence, cross entropy
- Functional derivatives
- Practice: MLE vs Maximum entropy

# High dimensional distributions

Probabilistic graphical models (PGMs) are about representing probability distributions over random variables

$$p(X) \equiv p(X_1, \dots, X_n)$$

$$2^n - 1$$

Assume  $x_i \in \{0, 1\}^n$

Naively, since there are  $2^n$  possible assignments to  $X_1, \dots, X_n$ , can represent this distribution completely using  $2^n - 1$  numbers, but quickly becomes intractable for large  $n$

PGMs are methods to represent these distributions more compactly, by exploiting *conditional independence*

CRF

## Recap: Bayesian networks (directed PGMs)

- A **Bayesian network** is specified by a directed acyclic graph  $G = (V, E)$  with:
  - ① One node  $i \in V$  for each random variable  $X_i$
  - ② One conditional probability distribution (CPD) per node,  $p(x_i | \mathbf{x}_{\text{Pa}(i)})$ , specifying the variable's probability conditioned on its parents' values

- Corresponds 1-1 with a particular factorization of the joint distribution:

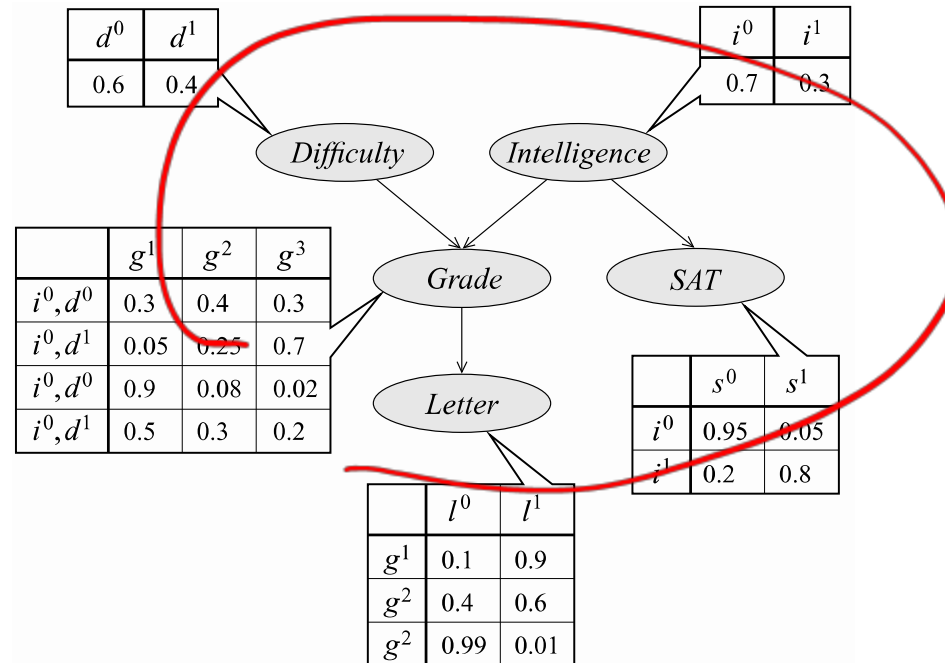
$$\underbrace{p(x_1, \dots, x_n)}_{2^n - 1} = \prod_{i \in V} p(x_i | \mathbf{x}_{\text{Pa}(i)}) = \prod_{i \in V} p(x_i)$$

*(Handwritten notes:  $2^n - 1$  above the first term,  $n$  to the right of the second term)*

- Powerful framework for designing *algorithms* to perform probability computations

# Recap: Example

- Consider the following Bayesian network:



- What is its joint distribution?

$$p(x_1, \dots, x_n) = \prod_{i \in V} p(x_i \mid \mathbf{x}_{\text{Pa}(i)})$$

$$p(d, i, g, s, l) = p(d)p(i)p(g \mid i, d)p(s \mid i)p(l \mid g)$$

# Entropy

# Entropy

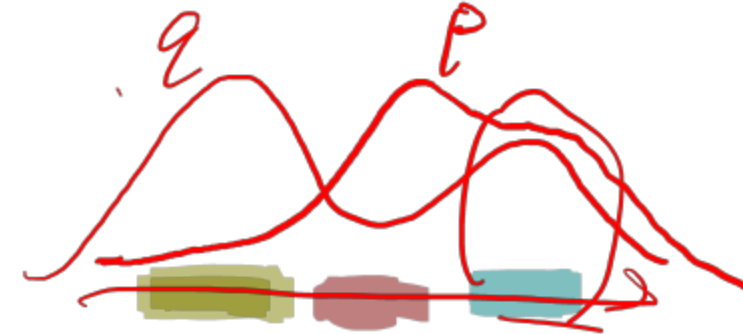
- Shannon entropy  $H(p) = - \sum_x p(x) \log p(x)$ 
  - The average level of "information", "surprise", or "uncertainty" inherent to the variable  $x$ 's possible outcomes

# KL Divergence

- Kullback-Leibler (KL) divergence: measures the closeness of two distributions  $p(x)$  and  $q(x)$

$$KL(q(x) || p(x)) = \sum_x q(x) \log \frac{q(x)}{p(x)}$$

*distance*



- a.k.a. Relative entropy
- $KL \geq 0$  (Jensen's inequality)

## Questions:

- If  $q$  is high and  $p$  is high in a region, then KL divergence is \_\_\_\_\_ in this region.
- If  $q$  is high and  $p$  is low in a region, then KL divergence is \_\_\_\_\_ in this region.
- If  $q$  is low in a region, then KL divergence is \_\_\_\_\_ in this region.

*low*  
*high*  
*low*



# KL Divergence

- Kullback-Leibler (KL) divergence: measures the closeness of two distributions  $p(x)$  and  $q(x)$

$$KL(q(x) || p(x)) = \sum_x q(x) \log \frac{q(x)}{p(x)}$$

- a.k.a. Relative entropy
- $KL \geq 0$  (Jensen's inequality)
- Intuitively:

- If  $q$  is high and  $p$  is high, then we are happy (i.e. low KL divergence)
- If  $q$  is high and  $p$  is low then we pay a price (i.e. high KL divergence).
- If  $q$  is low then we don't care (i.e. also low KL divergence, regardless of  $p$ )

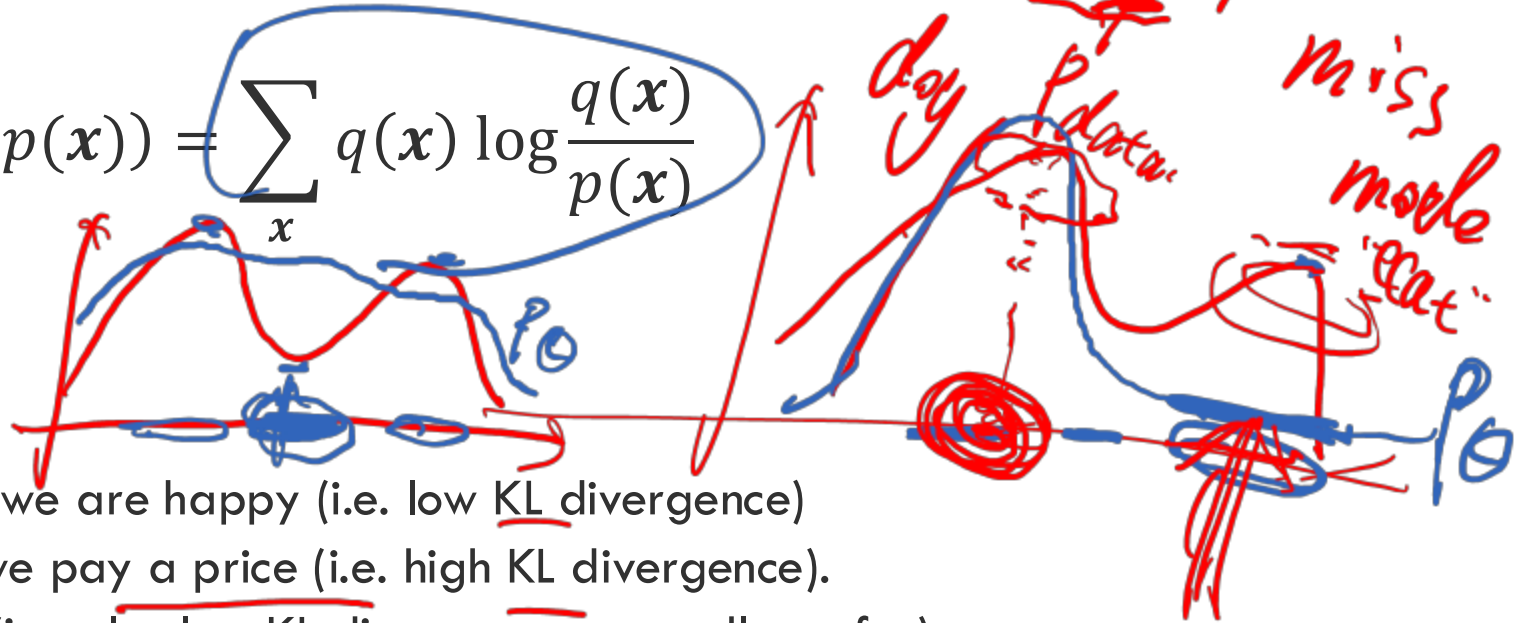
- not a true "distance":

- not commutative (symmetric)  $KL(p||q) \neq KL(q||p)$
- doesn't satisfy triangle inequality

①

GAN

$$\min_{\theta} KL(P_{\theta} || P_{data})$$



②

VAE

mode seeking

$$\min_{\theta} KL(P_{data} || P_{\theta})$$

# KL Divergence

- Kullback-Leibler (KL) divergence: measures the closeness of two distributions  $p(\mathbf{x})$  and  $q(\mathbf{x})$

$$\text{KL}(q(\mathbf{x}) \parallel p(\mathbf{x})) = \sum_{\mathbf{x}} q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})}$$

- a.k.a. Relative entropy

- Maximum likelihood estimation (MLE)  $\min_{\theta} - \mathbb{E}_{x \sim \tilde{p}(x)} \left[ \log p_{\theta}(x) \right]$

- **Question:** Show that MLE is minimizing the KL divergence between the empirical data distribution and the model distribution

# KL Divergence

$$\sum_x \tilde{p}(x) \log \frac{p(x)}{p_\theta(x)} = \sum_x \tilde{p}(x) \log \tilde{p}(x) - \sum_x \tilde{p}(x) \log p_\theta(x)$$

- Kullback-Leibler (KL) divergence: measures the closeness of two distributions  $p(x)$  and  $q(x)$

$$\text{KL}(q(x) \parallel p(x)) = \sum_x q(x) \log \frac{q(x)}{p(x)}$$

- a.k.a. Relative entropy

- Maximum likelihood estimation (MLE)

$$\min_{\theta} - \mathbb{E}_{x \sim \tilde{p}(x)} \left[ \log p_{\theta}(x) \right]$$

- **Question:** Show that MLE is minimizing the KL divergence between the empirical data distribution and the model distribution

$$\min_{\theta} \text{KL}(\tilde{p}(x) \parallel p_{\theta}(x)) = \underbrace{-\mathbb{E}_{\tilde{p}(x)} [\log p_{\theta}(x)]}_{\text{Cross entropy}} - \cancel{H(\tilde{p}(x))}$$

# Key Takeaways

- Probability  $p(\mathbf{x})$

- Bayes' rule 
$$p(\mathbf{y}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{y})p(\mathbf{y})}{p(\mathbf{x})}$$
  - prior, posterior

- Exponential family:
  - Gaussian, multinomial, categorical, ...

- Probabilistic graphical models: Bayesian networks

- KL Divergence
  - relation to Cross-entropy

$$\text{KL}(q(\mathbf{x}) || p(\mathbf{x})) = \sum_{\mathbf{x}} q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})}$$

# Functional Derivatives (Optional)

# Functional derivative

- $\nabla_q - \mathbb{H}(q) = \log q + 1$

max entropy

- Functional  $F(y)$ : an operator that takes a function  $y(x)$  and returns an output value  $F$

ER

x

- Functional derivative (aka, variational derivative): relates a change in a Functional  $F(y)$  to a change in the function  $y$

# Functional derivative

- Recall the conventional derivative  $\frac{dy}{dx}$ 
  - How much a function  $y(x)$  changes when we make a small change  $\Delta x$  to the input  $x$
  - Taylor expansion

$$y(x + \epsilon) = y(x) + \frac{dy}{dx} \epsilon + O(\epsilon^2)$$

- Functional derivative
  - How much a functional  $F[y]$  changes when we make a small change  $\epsilon \eta(x)$  to the function  $y(x)$

$$F[y(x) + \epsilon \eta(x)] = F[y(x)] + \epsilon \int \frac{\delta F}{\delta y(x)} \eta(x) dx + O(\epsilon^2)$$

- A function  $y(x)$  that maximizes (or minimizes) a functional  $F[y]$  must satisfy

$$\frac{\delta F}{\delta y(x)} = 0 \text{ for all } x$$

# Functional derivative

$$F[y(x) + \epsilon\eta(x)] = F[y(x)] + \epsilon \int \frac{\delta F}{\delta y(x)} \eta(x) dx + O(\epsilon^2)$$

- Consider a functional that is defined by an integral over a function  $G(y, x)$

$$F[y] = \int G(y, x) dx$$

- Ex.1,  $-H(q) = \int q(x) \log q(x) dx$

- $G = q(x) \log q(x)$

- Consider variations in the function  $y(x)$ ,

$$H(q) = \int q(x) \log q(x) dx$$

$G(q, x)$

$$F[y + \epsilon\eta(x)] = F[y(x)] + \epsilon \int \frac{\partial G}{\partial y} \eta(x) dx + O(\epsilon^2)$$

$$\frac{\partial q \cdot \log q}{\partial q} = \log q + q \cdot \frac{1}{q} = \log q + 1 \quad \frac{\partial H(q)}{\partial q} = \frac{\partial G}{\partial q}$$



# Functional derivative

$$F[y(x) + \epsilon\eta(x)] = F[y(x)] + \epsilon \int \frac{\delta F}{\delta y(x)} \eta(x) dx + O(\epsilon^2)$$

- Consider a functional that is defined by an integral over a function  $G(y, x)$

- Ex.1,  $-\mathbb{H}(q) = \int q(x) \log q(x) dx$ 
  - $G = q(x) \log q(x)$

- Consider variations in the function  $y(x)$ ,

$$F[y + \epsilon\eta(x)] = F[y(x)] + \epsilon \int \frac{\partial G}{\partial y} \eta(x) dx + O(\epsilon^2)$$

# Practice: Maximum likelihood vs Maximum Entropy (Optional)

# Supervised Maximum Likelihood

- Model to be learned  $p_{\theta}(x)$   *$P_{\theta}$  (image, label)*
- Observe full data  $\mathcal{D} = \{x^*\}$ 
  - i.i.d: independent, identically distributed
- Maximum Likelihood Estimation (MLE)
  - The most classical learning algorithm

$$\min_{\theta} - \mathbb{E}_{x^* \sim \mathcal{D}} \left[ \log p_{\theta}(x^*) \right]$$

- MLE is closely connected to the Maximum Entropy (MaxEnt) principle

# Recap: Exponential Family

- A distribution

$$p_{\theta}(\mathbf{x}) = h(\mathbf{x}) \exp\{\theta \cdot T(\mathbf{x})\} / Z(\theta)$$

is an exponential family distribution

- $\theta \in R^d$ : natural (canonical) parameter
  - $T(\mathbf{x}) \in R^d$ : sufficient statistics, features of data  $\mathbf{x}$
  - $Z(\theta) = \sum_{\mathbf{x}, y} h(\mathbf{x}) \exp\{\theta \cdot T(\mathbf{x})\}$ : normalization factor
- Examples: Bernoulli, multinomial, Gaussian, Poisson, gamma,...

# Maximum Likelihood for Exponential Family

$m(\mathbf{x})$ : the number of times  $\mathbf{x}$  is observed in  $D$

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}; \mathcal{D}) &= \sum_{\mathbf{x}} m(\mathbf{x}) \log p(\mathbf{x} | \boldsymbol{\theta}) \\ &= \sum_{\mathbf{x}} m(\mathbf{x}) \left( \sum_i \theta_i T_i(\mathbf{x}) - \log Z(\boldsymbol{\theta}) \right) \\ &= \sum_{\mathbf{x}} m(\mathbf{x}) \sum_i \theta_i T_i(\mathbf{x}) - N \log Z(\boldsymbol{\theta})\end{aligned}$$

- Take gradient and set to 0

$$\begin{aligned}\frac{\partial}{\partial \theta_i} \mathcal{L}(\boldsymbol{\theta}; \mathcal{D}) &= \sum_{\mathbf{x}} m(\mathbf{x}) T_i(\mathbf{x}) - N \frac{\partial}{\partial \theta_i} \log Z(\boldsymbol{\theta}) \\ &= \sum_{\mathbf{x}} m(\mathbf{x}) T_i(\mathbf{x}) - N \sum_{\mathbf{x}} p(\mathbf{x} | \boldsymbol{\theta}) T_i(\mathbf{x})\end{aligned}$$

$$\Rightarrow \sum_{\mathbf{x}} p(\mathbf{x} | \boldsymbol{\theta}) T_i(\mathbf{x}) = \sum_{\mathbf{x}} \frac{m(\mathbf{x})}{N} T_i(\mathbf{x}) = \sum_{\mathbf{x}} \tilde{p}(\mathbf{x} | \boldsymbol{\theta}) T_i(\mathbf{x})$$

At MLE, the expectations of the sufficient statistics under the model must match empirical feature average

# Maximum Entropy (MaxEnt)

- Given  $\mathcal{D}$ , to estimate  $p(\mathbf{x})$
- We can approach the problem from an entirely different point of view. Begin with some fixed feature expectations:

$$\sum_{\mathbf{x}} p(\mathbf{x}) T_i(\mathbf{x}) = \sum_{\mathbf{x}} \frac{m(\mathbf{x})}{N} T_i(\mathbf{x}) := \alpha_i$$

- There may exist many distributions which satisfy them. Which one should we select?
  - MaxEnt principle: the most uncertain or flexible one, i.e., the one with maximum entropy
- This yields a new optimization problem:
  - This is a variational definition of a distribution!

~~\_\_\_\_\_~~  
Stem is defined as a solution  
to a opt problem

$$\begin{aligned} \max_p \quad & \underline{H(p(\mathbf{x}))} \\ \text{s.t.} \quad & \left\{ \begin{aligned} \sum_{\mathbf{x}} p(\mathbf{x}) T_i(\mathbf{x}) &= \alpha_i \\ \sum_{\mathbf{x}} p(\mathbf{x}) &= 1 \end{aligned} \right. \end{aligned}$$

# Solution to the MaxEnt Problem

- To solve the MaxEnt problem, we use Lagrange multipliers:

$$\max_{\theta, \mu} \min_{p(x)} L = - \sum_{\mathbf{x}} p(\mathbf{x}) \log p(\mathbf{x}) - \sum_i \theta_i \left( \sum_{\mathbf{x}} p(\mathbf{x}) T_i(\mathbf{x}) - \alpha_i \right) - \mu \left( \sum_{\mathbf{x}} p(\mathbf{x}) - 1 \right)$$

$$L = \int a(p, x) dx$$

$$\partial \sum_i \theta_i p(x) T_i(x)$$

## Solution to the MaxEnt Problem

- To solve the MaxEnt problem, we use Lagrange multipliers:

$$\max_{\theta, \mu} \min_{p(x)} L = - \sum_{\mathbf{x}} p(\mathbf{x}) \log p(\mathbf{x}) - \sum_i \theta_i \left( \sum_{\mathbf{x}} p(\mathbf{x}) T_i(\mathbf{x}) - \alpha_i \right) - \mu \left( \sum_{\mathbf{x}} p(\mathbf{x}) - 1 \right)$$

$$\frac{\partial L}{\partial p(\mathbf{x})} = 1 + \log p(\mathbf{x}) - \sum_i \theta_i T_i(\mathbf{x}) - \mu = 0$$

$$p^*(\mathbf{x}) = e^{\mu-1} \exp \left\{ \sum_i \theta_i f_i(\mathbf{x}) \right\}$$

$$Z(\boldsymbol{\theta}) = e^{\mu-1} = \sum_{\mathbf{x}} \exp \left\{ \sum_i \theta_i f_i(\mathbf{x}) \right\} \quad \left( \text{since } \sum_{\mathbf{x}} p^*(\mathbf{x}) = 1 \right)$$

$$p(\mathbf{x} | \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \exp \left\{ \sum_i \theta_i T_i(\mathbf{x}) \right\}$$

$\partial p$

$\frac{\partial \mu p(x)}{\partial p(x)}$



# Solution to the MaxEnt Problem

- To solve the MaxEnt problem, we use Lagrange multipliers:

$$\max_{\theta, \mu} \min_{p(\mathbf{x})} L = - \sum_{\mathbf{x}} p(\mathbf{x}) \log p(\mathbf{x}) - \sum_i \theta_i \left( \sum_{\mathbf{x}} p(\mathbf{x}) T_i(\mathbf{x}) - \alpha_i \right) - \mu \left( \sum_{\mathbf{x}} p(\mathbf{x}) - 1 \right)$$

$$\frac{\partial L}{\partial p(\mathbf{x})} = 1 + \log p(\mathbf{x}) - \sum_i \theta_i T_i(\mathbf{x}) - \mu$$

$$p^*(\mathbf{x}) = e^{\mu-1} \exp \left\{ \sum_i \theta_i f_i(\mathbf{x}) \right\}$$

$$Z(\boldsymbol{\theta}) = e^{\mu-1} = \sum_{\mathbf{x}} \exp \left\{ \sum_i \theta_i f_i(\mathbf{x}) \right\} \quad \left( \text{since } \sum_{\mathbf{x}} p^*(\mathbf{x}) = 1 \right)$$

$$p(\mathbf{x} | \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \exp \left\{ \sum_i \theta_i T_i(\mathbf{x}) \right\}$$

- So feature constraints + MaxEnt  $\Rightarrow$  exponential family.

- Problem is strictly convex w.r.t.  $p(\mathbf{x})$ , so solution is unique.

# Solution to the MaxEnt Problem

- To solve the MaxEnt problem, we use Lagrange multipliers:

$$\max_{\theta, \mu} \min_{p(\mathbf{x})} L = - \sum_{\mathbf{x}} p(\mathbf{x}) \log p(\mathbf{x}) - \sum_i \theta_i \left( \sum_{\mathbf{x}} p(\mathbf{x}) T_i(\mathbf{x}) - \alpha_i \right) - \mu \left( \sum_{\mathbf{x}} p(\mathbf{x}) - 1 \right)$$

$$p(\mathbf{x} | \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \exp \left\{ \sum_i \theta_i T_i(\mathbf{x}) \right\}$$

plug  $p(\mathbf{x} | \boldsymbol{\theta})$  back into  $L$ , and since  $\sum_{\mathbf{x}} \frac{m(\mathbf{x})}{N} T_i(\mathbf{x}) := \alpha_i$ :

$$\max_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) = \sum_{\mathbf{x}} m(\mathbf{x}) \sum_i \theta_i T_i(\mathbf{x}) - N \log Z(\boldsymbol{\theta})$$

- Recovers precisely the MLE problem of exponential family

- So feature constraints + MaxEnt  $\Rightarrow$  **exponential family**.
- Problem is strictly convex w.r.t.  $p(\mathbf{x})$ , so solution is unique.

# Constraints from Data

- We have seen a case of **convex duality**:
  - In one case, we assume exponential family and show that Maximum Likelihood implies model expectations must match empirical expectations.
  - In the other case, we assume model expectations must match empirical feature counts and show that MaxEnt implies exponential family distribution.

## A more general MaxEnt problem

$$\min_p \text{KL}(p(\mathbf{x}) \| h(\mathbf{x}))$$

$$\stackrel{\text{def}}{=} \sum_{\mathbf{x}} p(\mathbf{x}) \log \frac{p(\mathbf{x})}{h(\mathbf{x})} = -\text{H}(p) - \sum_{\mathbf{x}} p(\mathbf{x}) \log h(\mathbf{x})$$

$$\text{s.t.} \quad \sum_{\mathbf{x}} p(\mathbf{x}) T_i(\mathbf{x}) = \alpha_i$$

$$\sum_{\mathbf{x}} p(\mathbf{x}) = 1$$

$$\Rightarrow p(\mathbf{x} | \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} h(\mathbf{x}) \exp \left\{ \sum_i \theta_i T_i(\mathbf{x}) \right\}$$

# Summary

- Maximum entropy is dual to maximum likelihood of exponential family distributions
- This provides an alternative view of the problem of fitting a model into data:
  - The data instances in the training set are treated as constraints, and the learning problem is treated as a constrained optimization problem.
  - We'll revisit this optimization-theoretic view of learning repeatedly in the future!

$$\begin{aligned} \max_p \quad & H(p(\mathbf{x})) = - \sum_{\mathbf{x}} p(\mathbf{x}) \log p(\mathbf{x}) \\ \text{s.t.} \quad & \sum_{\mathbf{x}} p(\mathbf{x}) T_i(\mathbf{x}) = \alpha_i \\ & \sum_{\mathbf{x}} p(\mathbf{x}) = 1 \end{aligned}$$

$r(x, y) \geq 0$

# Key Takeaways

- Probability
  - Bayes' rule
  - Exponential family
  - Probabilistic graphical models: Bayesian networks
  - KL divergence
- Functional derivative
- Convex duality between MLE and MaxEnt (optional)

**Questions?**