

Machine Learning 2

Conditional Independence

Patrick Forré

Recap: Marginals & Conditionals

- Let X, Y, Z be random variables with joint distribution $p(x, y, z)$.
- Marginal distribution of Z :

$$p(z) = \int p(x, y, z) dx dy$$

- (Marginal) joint distribution of X and Y given $Z = z$:

$$p(x, y | z) = \begin{cases} \frac{p(x, y, z)}{p(z)}, & p(z) > 0 \\ 0, & p(z) = 0 \end{cases}$$

Recap: Independence

- Let X, Y be random variables with joint distribution $p(x, y)$.
We say that X is **independent** of Y iff:

$$\forall x, y, \cancel{z}: p(x, y) = p(x) \cdot p(y)$$

- In symbols:

$$X \perp\!\!\!\perp Y$$

Conditional Independence - Definition

- Let X, Y, Z be random variables with joint distribution $p(x, y, z)$. We say that X is independent of Y conditioned on Z iff:

$$\forall x, y, z : p(x, y | z) = p(x|z) \cdot p(y|z)$$

(actually, we only need to check that for z with $p(z) > 0$).

- In symbols:

$$X \perp\!\!\!\perp Y | Z$$

Conditional Independence - Remarks

- We consider (unconditional) independence as a special case of conditional independence (as we need to say what we condition on anyways):

$$X \perp\!\!\!\perp Y \iff X \perp\!\!\!\perp Y | \emptyset \iff X \perp\!\!\!\perp Y | (\text{const. RV})$$

- Independence heuristically means that knowing Y does not give us any information about X .
- Conditional Independence of X and Y given Z means that there is no additional information in the state of Y beyond the information already in Z to determine the state of X .

Conditional Independence - Example 1

- Alice and Bob each independently throw a coin.
- They try to guess what the other person has thrown.
- Is that possible by just looking at their own coin?

A $\perp\!\!\!\perp$ B

- Eve just happens to see both coin outcomes and tells Bob: “I don’t tell you what Alice has, the outcomes are different anyways.”
- What about now?

A $\cancel{\perp\!\!\!\perp}$ B | E



Conditional Independence - Example 2

- Alice, Bob and Casper each throw a coin.
- Alice tries to guess the number of ‘tails’ T .
- Alice has a ‘head’, what can she tell?



A ~~H~~ T

- Eve just saw that Alice and Bob have 2 ‘heads’ together.
- Can Alice now tell if there now 2 or 3 ‘tails’ in total?

A ~~H~~ T | E

Exercise: Separoid Axioms

- Prove the following claims for random variables X, Y, Z, W :

- Redundancy: $X \perp\!\!\!\perp Y | X$ always holds.

- Symmetry: $X \perp\!\!\!\perp Y | Z \iff Y \perp\!\!\!\perp X | Z$

- Decomposition: $X \perp\!\!\!\perp (Y, W) | Z \implies X \perp\!\!\!\perp Y | Z$

- Weak Union: $X \perp\!\!\!\perp (Y, W) | Z \implies X \perp\!\!\!\perp W | (Y, Z)$

- Contraction: $X \perp\!\!\!\perp W | (Y, Z) \wedge X \perp\!\!\!\perp Y | Z \implies X \perp\!\!\!\perp (Y, W) | Z$

- So we have the equivalence:

$$X \perp\!\!\!\perp (Y, W) | Z \iff X \perp\!\!\!\perp W | (Y, Z) \wedge X \perp\!\!\!\perp Y | Z$$

Summary

- Conditional Independence: $X \perp\!\!\!\perp Y|Z$ defined to hold if:

$$\forall x, y, z : p(x, y | z) = p(x | z) \cdot p(y | z)$$

- Example 1 shows: $X \perp\!\!\!\perp Y \not\Rightarrow X \perp\!\!\!\perp Y|Z$
- Example 2 shows: $X \perp\!\!\!\perp Y|Z \not\Rightarrow X \perp\!\!\!\perp Y$

Machine Learning 2

Information Theory
- Entropy

Patrick Forré

Let's Play: "Who am I?"

- Find out who you represent in a minimal number of YES/NO-questions.
- What is the minimal number of YES/NO-questions required averaged over all possible characters?

→ Entropy



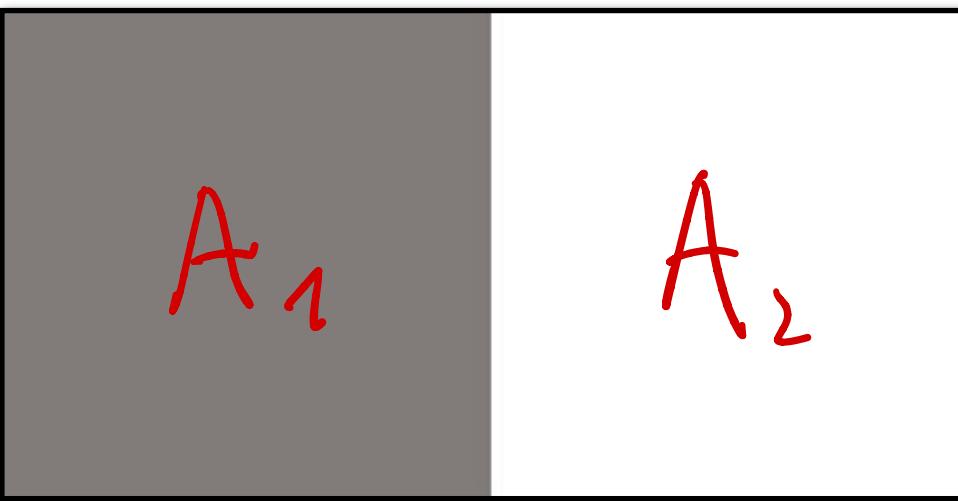
- need to know whole space of possible characters
(with frequencies)

Entropy - Informal Definition

- Let X be a random variable.
- Entropy of $X :=$ Minimal bits (=YES/NO-questions) needed to encode state X averaged over all possible states that X can take weighted their probability/frequency.
- $H(X) := \mathbb{E}[\#\text{bits}(X)]$

Entropy - Examples

- Alice selects field.
- How many YES/NO-questions does Bob need to ask to find that field?



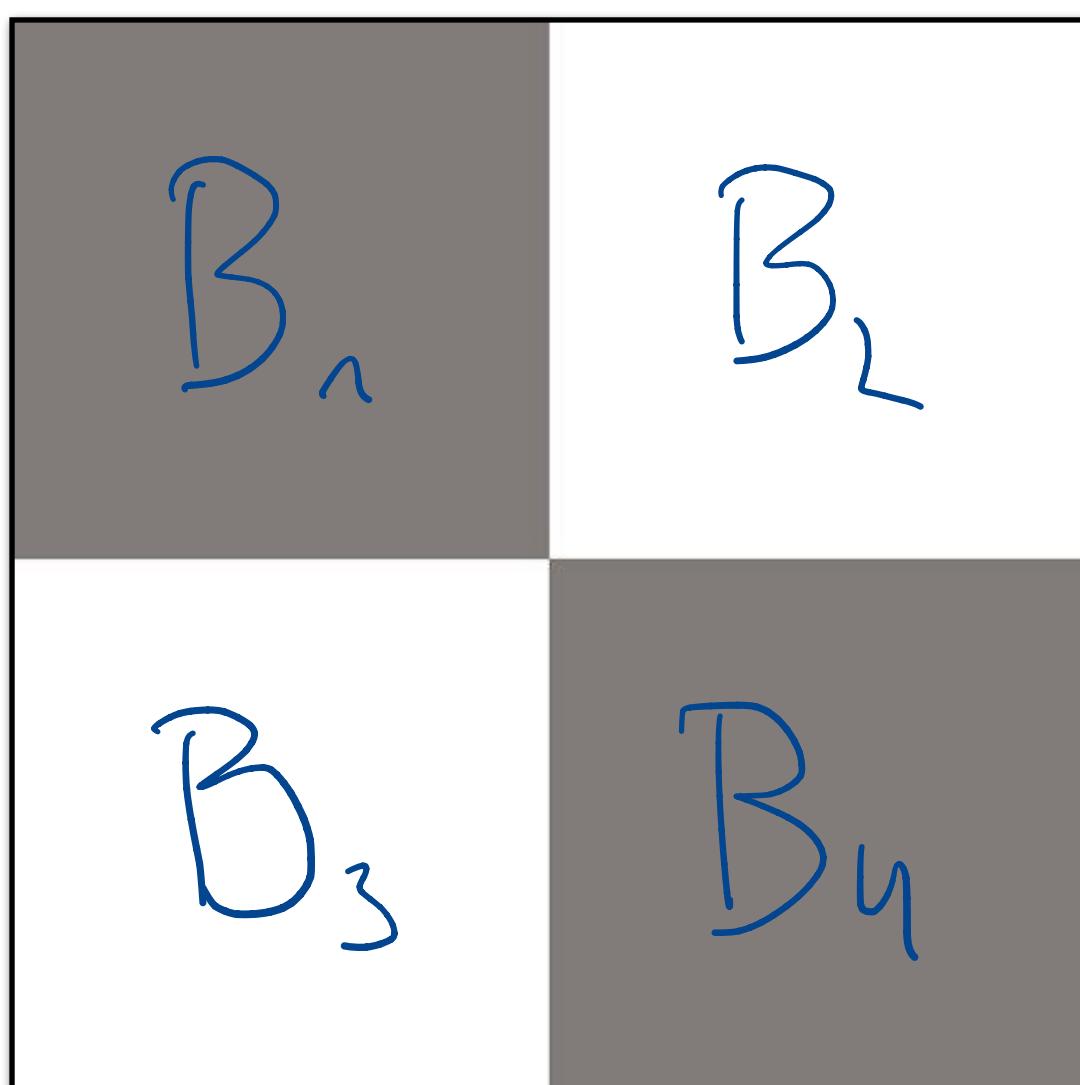
$$\# \text{bits}(A) = 1$$

$$p(A_1) = \frac{1}{2}$$

$$\# \text{bits}(A_2) = 1$$

$$p(A_2) = \frac{1}{2}$$

$$H(A) = 1$$



$$\# \text{bits}(B_i) = 2$$

$$p(B_i) = \frac{1}{4} = \frac{1}{2^2}$$

$$\begin{aligned} H(B) &= \mathbb{E}[\# \text{bits}(B)] \\ &= \sum_{i=1}^4 p(B_i) \cdot \# \text{bits}(B_i) \\ &= \sum_{i=1}^4 p(B_i) \cdot 2 \end{aligned}$$

$$= 2$$

Entropy - Examples

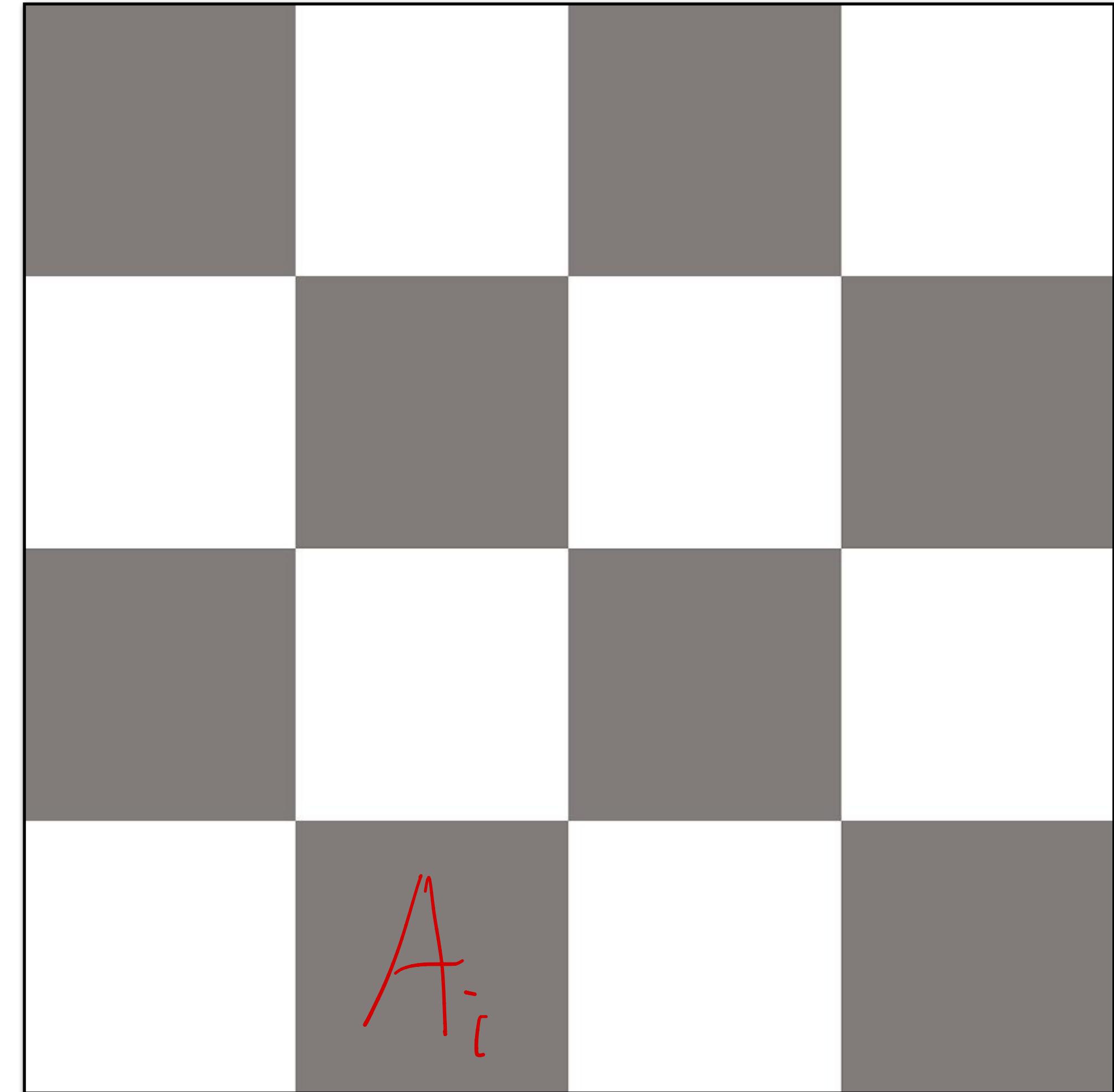
- Alice selects field.
- How many YES/NO-questions does Bob need to ask to find that field?

$$\#\text{bits}(A_i) = 4 \quad , \quad H(A) = 4$$

$$P(A_i) = \frac{1}{2^4}$$

$$P(A_i) = \frac{1}{2^{\#\text{bits}}}$$

$$\sim \boxed{\#\text{bits}(A_i) = -\log_2 P(A_i)}$$



Entropy - Examples

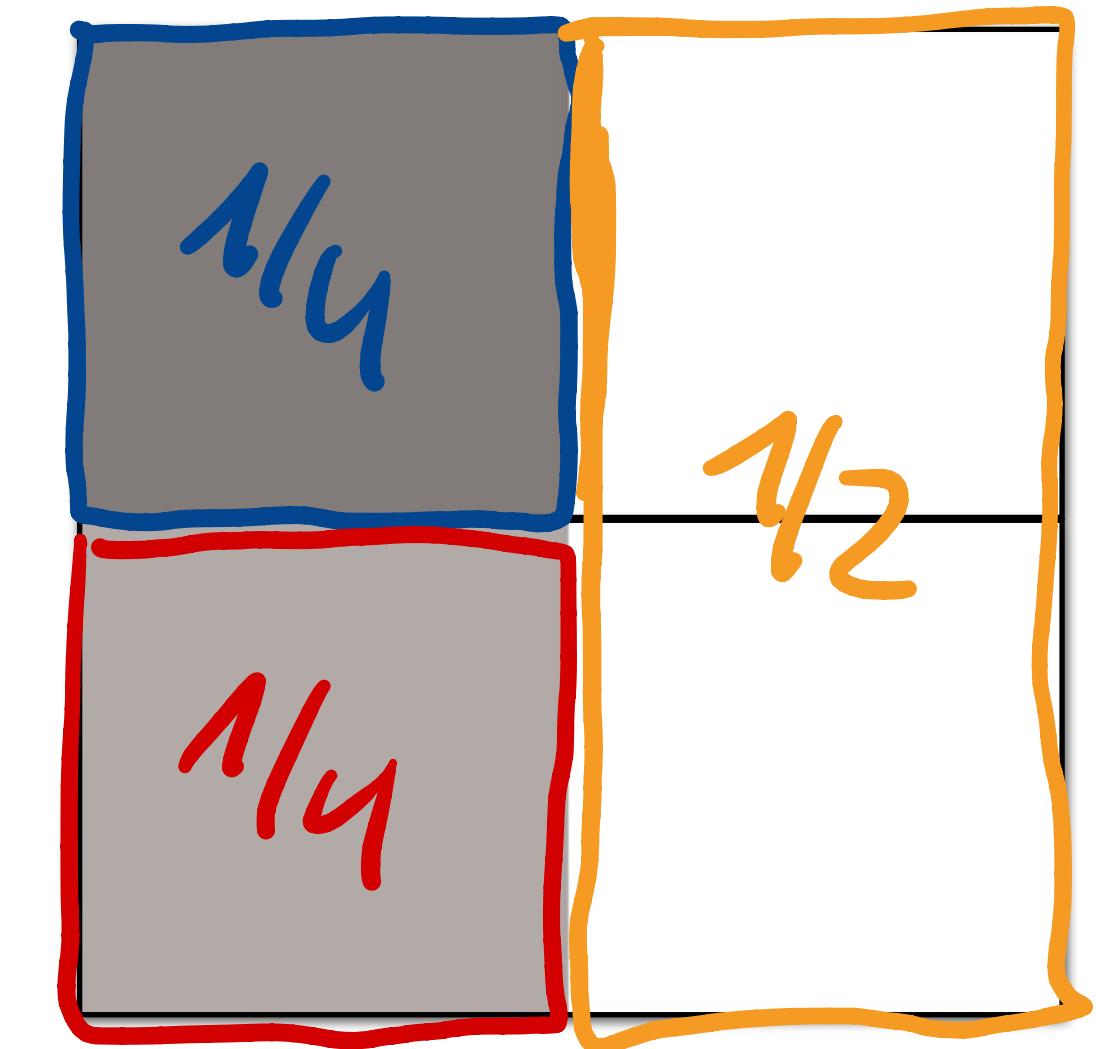
- Alice selects field.
- How many YES/NO-questions does Bob need to ask to find that field?

$$\# \text{bits}(A_1) = 2, \quad p(A_1) = \frac{1}{4} = \frac{1}{2^{\# \text{bits}}}$$

$$\# \text{bits}(A_2) = 2, \quad p(A_2) = \frac{1}{4} = \frac{1}{2^{\# \text{bits}}}$$

$$\# \text{bits}(A_3) = 1, \quad p(A_3) = \frac{1}{2} = \frac{1}{2^{\# \text{bits}}}$$

$$H(A) = \sum_{i=1}^3 p(A_i) \cdot \# \text{bits}(A_i) = \frac{1}{4} \cdot 2 + \frac{1}{4} \cdot 2 + \frac{1}{2} \cdot 1 \\ = \frac{3}{2} \quad \epsilon [1, 2]$$



Entropy - Examples

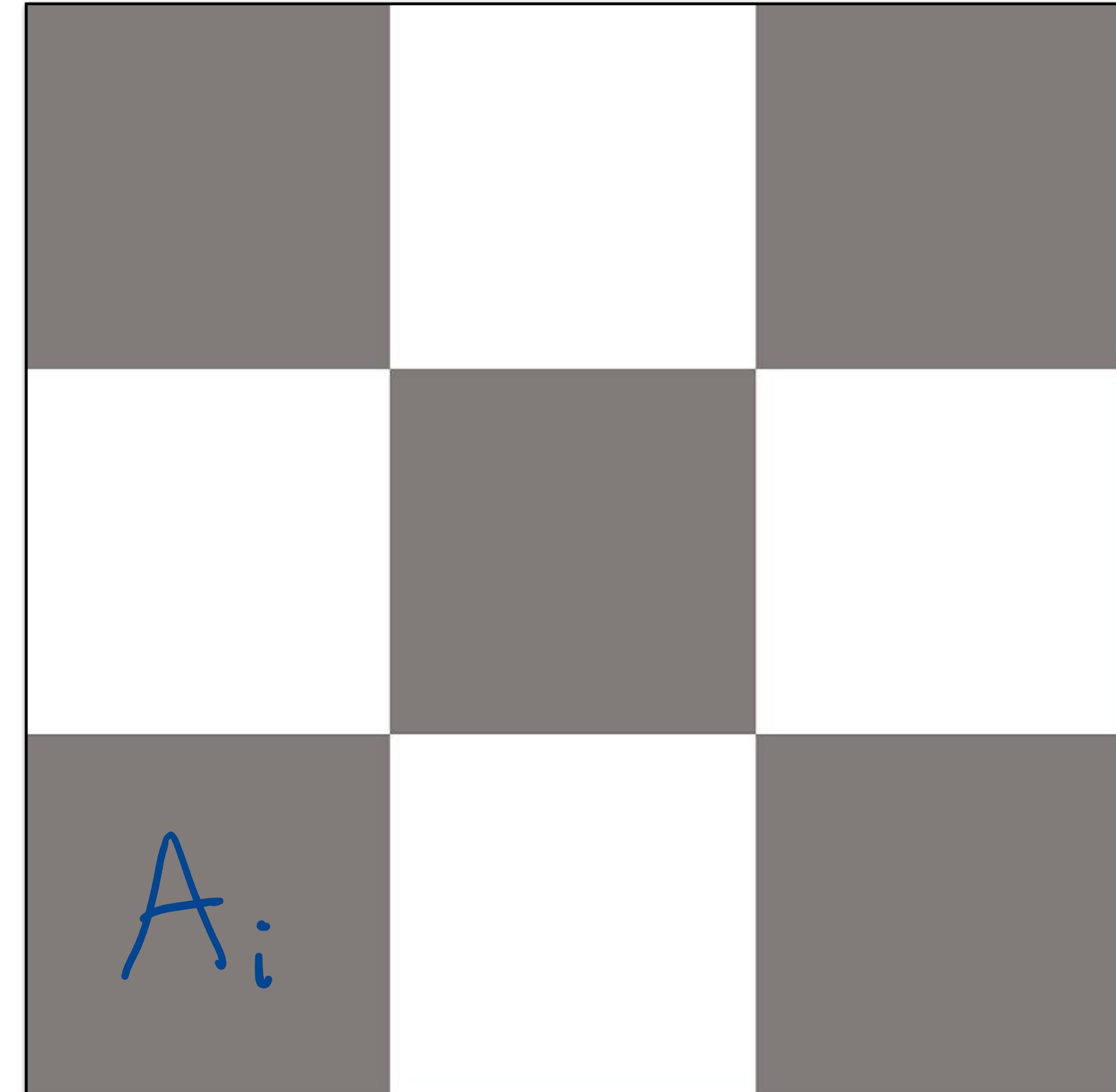
- Alice selects field.
- How many YES/NO-questions does Bob need to ask to find that field?

$$p(A_i) = \frac{1}{9}$$

$$\#\text{bits}(A_i) = -\log_2 p(A_i)$$

Fano-Shannon formula / code

$$\begin{aligned}H(A) &= \mathbb{E}[\#\text{bits}(A)] = \mathbb{E}[-\log p(A)] \\&= -\log_2 \frac{1}{9} = \log_2 9 \approx 3.17\end{aligned}$$



Entropy - Definition

- Let X be a random variable with distribution $p(x)$.

Shannon Entropy
(X discrete)

- Entropy of X :

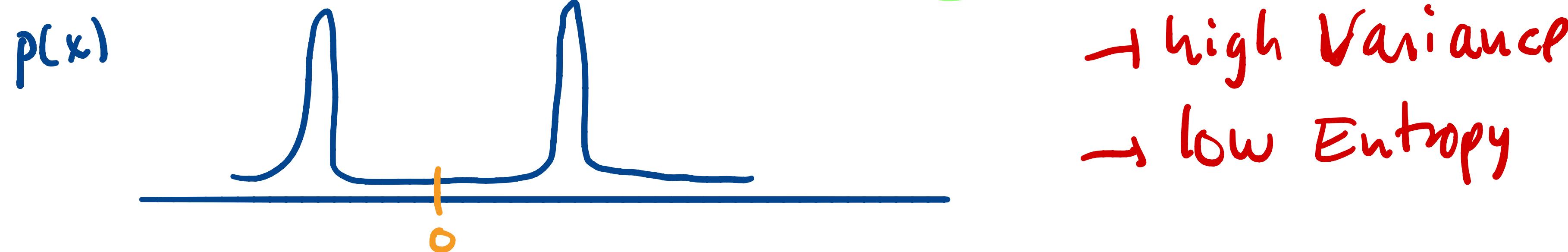
$$H(X) := \mathbb{E}[-\log p(x)]$$

Differential Entropy
(X continuous)

- The choice of the basis of log does not really matter, but one should be consistent with it. Usually one uses basis 2 to measure ‘bits’ or basis e to measure ‘nats’.

Entropy - Remarks

- The Entropy measures spread/uncertainty in the distribution similar to the Variance, but does not need a center of mass $\mathbb{E}[X]$ as a reference point like the Variance does.
- Entropy is also meaningful in the multimodal case.



- $H(X)$ small $\rightarrow X$ concentrated in small area.
- $H(X)$ big $\rightarrow X$ is scattered around.
- $H(X) \geq 0$ for discrete variables, but not in general.

Entropy - Exponential Families

- $p(x|\eta) = h(x) \cdot \exp(\eta^\top T(x) - A(\eta))$

$$H(X) = \mathbb{E}_{\eta}[-\log p(x|\eta)] = \mathbb{E}\left[-\log h(x) - (\eta^\top T(x) - A(\eta))\right]$$

$$= -\mathbb{E}[\log h(x)] - \eta^\top \underbrace{\mathbb{E}[T(x)]}_{\nabla_{\eta} A(\eta)} + A(\eta)$$

$$= A(\eta) - \eta^\top \cdot \nabla_{\eta} A(\eta) - \mathbb{E}[\log h(x)]$$

(in nats)

Entropy - D-dimensional Gaussians

- $p(x | \mu, \Sigma) = \det(2\pi\Sigma)^{-\frac{1}{2}} \cdot \exp\left(-\frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu)\right)$

$$\begin{aligned} H(X) &= \mathbb{E}[-\log p(x | \mu, \Sigma)] \\ &= \frac{1}{2} \log \det(2\pi\Sigma) + \frac{1}{2} \mathbb{E}\left[\underbrace{(x - \mu)^\top \Sigma^{-1} (x - \mu)}_{\sum \text{Tr}[\Sigma^{-1} (x - \mu) (x - \mu)^\top]} \right] \cdot \log e \\ &= \frac{1}{2} \log \det(2\pi\Sigma) + \frac{1}{2} \text{Tr}\left[\Sigma^{-1} \mathbb{E}[(x - \mu) (x - \mu)^\top]\right] \cdot \log e \\ &= \frac{1}{2} \log \det(2\pi\Sigma) + \frac{1}{2} D \cdot \log e \\ &= \frac{1}{2} \log \det(2\pi \cdot e \cdot \Sigma) \end{aligned}$$

Machine Learning 2

**Information Theory
- Maximum Entropy Principle**

Patrick Forré

Maximum Entropy Principle

- “Whenever you don’t know what probability distribution to pick, then take a distribution expressing the highest amount of uncertainty (=Entropy), given all (known) constraints.”:

$$p^* \in \underset{p \in \mathcal{C}}{\operatorname{arg\,max}} H(p)$$

where $H(p) = \mathbb{E}[-\log p(X)]$ is the entropy for $X \sim p$ and \mathcal{C} the constraint set.

- Used for choosing a prior in Bayesian statistics.
- Used in classical statistical mechanics in physics.

Exercise: Deriving Exponential Families

- We want to derive a distribution $p(x)$ for X that satisfies the constraints:

$$\mathbb{E}_p[T_j(X)] = \tau_j \quad \text{for } j = 1, \dots, s \quad \text{for given functions } T_j \text{ and numbers } \tau_j.$$

- The Lagrangian for the maximum entropy principle is then:

$$\mathcal{L} = H(p) + \sum_{j=1}^s \eta_j \cdot \left(\mathbb{E}_p[T_j(X)] - \tau_j \right) + (1 - A) \cdot \left(\int p(x) dx - 1 \right)$$

with suggestively chosen Lagrange multipliers η_j , and ~~A~~ $(1 - A)$

- The solution of the maximum entropy principle is then of exponential family form with the suggested sufficient statistics and natural parameters, $h(x) = 1$, etc.

Machine Learning 2

**Information Theory
- Relative Entropy**

Patrick Forré

Recap: Jensen's Inequality

- Let X be a real-valued random variable and φ a convex function (e.g. if $\varphi''(x) \geq 0$ for all x). Then we have the inequality:

$$\varphi(\mathbb{E}[X]) \leq \mathbb{E}[\varphi(X)]$$

e.g. $|\mathbb{E}[X]| \leq \mathbb{E}|X|$

- In case φ is strictly convex (e.g. if $\varphi''(x) > 0$ for all x) then: equality holds in the above inequality if and only if X is (almost-surely) a constant random variable (with value $\mathbb{E}[X]$).

Kullback-Leibler Divergence

- Let $p(x)$ and $q(x)$ be two probability distribution on the same space. Then the **relative entropy** or **Kullback-Leibler divergence** of p and q is:

$$\text{KL}(p\|q) := \begin{cases} \int p(x) \cdot \log \frac{p(x)}{q(x)} dx & , \text{ continuous case} \\ \sum_x p(x) \log \frac{p(x)}{q(x)} & , \text{ discrete case} \\ \mathbb{E}_{x \sim p} [\log \frac{p(x)}{q(x)}] & \text{cross-entropy} \quad \text{entropy.} \end{cases}$$
$$= -\mathbb{E}_p[\log q] - (\mathbb{E}_p[-\log p]) = \overbrace{C(p||q)}^{} - H(p)$$

- Conventions: $0 \cdot \infty := 0$, $1/\infty := 0$, $1/0 := \infty$, $0/0 := 0$.

Fundamental Inequality of Information Theory

- Let $p(x)$ and $q(x)$ be two probability distribution on the same space. Then we always have the **inequality**:

$$\text{KL}(p||q) \geq 0$$

- Furthermore, **equality** holds if and only if $p(x) = q(x)$ for (p -almost) all x .
- Proof: $-\log$ is strictly convex. By Jensen's inequality we get:

$$\begin{aligned} \text{KL}(p||q) &= E_p[\log \frac{q}{p}] = E_p[-\log \frac{p}{q}] \stackrel{\text{Jensen}}{\geq} -\log E_p\left[\frac{q}{p}\right] = -\log \int p(x) \frac{q(x)}{p(x)} dx \\ &= -\log 1 = 0 \end{aligned}$$

Kullback-Leibler Divergence - Remarks

- KL is **not symmetric** in the arguments and can take value ∞ .
- KL measures the average **additional** bits needed to encode the location of x sampled from “true” distribution $p(x)$, but using “proposal” distribution $q(x)$ (instead of $p(x)$)
- Minimizing the KL divergence between some distributions can be described as the **principle of minimal relative entropy**.

Deriving Maximum Likelihood Estimation

- Let x_1, \dots, x_N be i.i.d. data with true distribution $q(x)$.
- Let $\{ p(x|\theta) | \theta \in \Theta \}$ be a statistical model aimed to describe $q(x)$.
- The principle of minimal relative entropy tells us to pick:

$$\theta^* \in \arg \min_{\theta \in \Theta} KL(q(x) \| p(x|\theta)) = \mathbb{E}_q \left[\log \frac{q(x)}{p(x|\theta)} \right]$$
$$\approx \frac{1}{N} \sum_{n=1}^N \log \frac{q(x_n)}{p(x_n|\theta)} = -\frac{1}{N} \sum_{n=1}^N \log p(x_n|\theta) + \text{const.}$$

Where we used the law of large numbers.

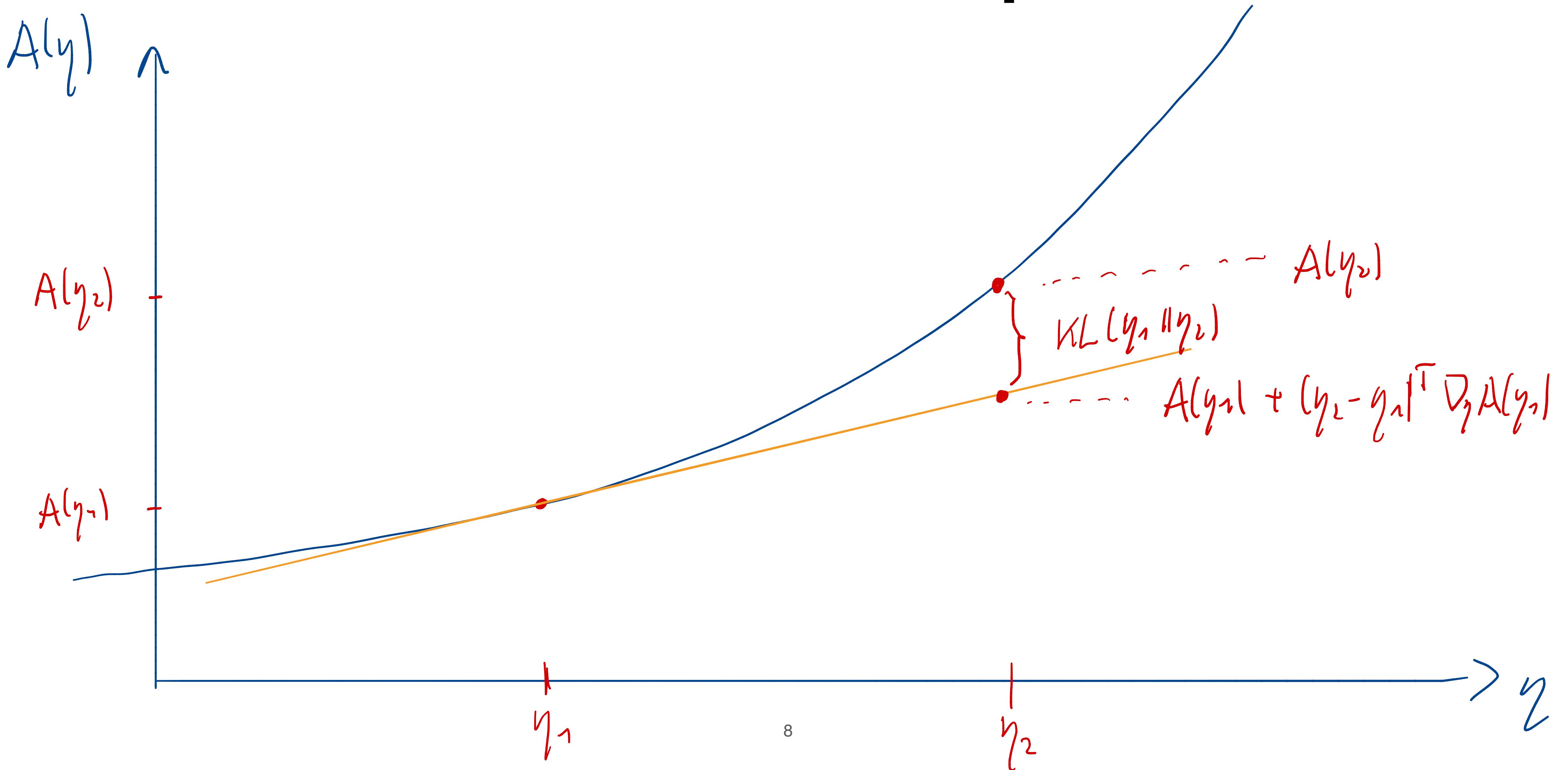
- Maximum likelihood estimation: $\hat{\theta} \in \arg \max_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N \log p(x_n|\theta)$

Example: Exponential Families

- Exponential family: $p(x|\eta) = h(x) \cdot \exp(\eta^\top T(x) - A(\eta))$
- Pick two natural parameters η_1, η_2 .

$$\begin{aligned} \text{KL}(p(X|\eta_1) \parallel p(X|\eta_2)) &= \mathbb{E}_{\eta_1} \left[\log \frac{p(x|\eta_1)}{p(x|\eta_2)} \right] \\ &= \mathbb{E}_{\eta_1} \left[(\eta_1 - \eta_2)^\top T(x) - A(\eta_1) + A(\eta_2) \right] = (\eta_1 - \eta_2)^\top \underbrace{\mathbb{E}_{\eta_1} [T(x)]}_{\nabla_{\eta_1} A(\eta_1)} - A(\eta_1) + A(\eta_2) \\ &= A(\eta_2) - [(\eta_2 - \eta_1)^\top \nabla_{\eta_1} A(\eta_1) + A(\eta_1)] \end{aligned}$$

Information Geometric Interpretation of KL



Machine Learning 2

**Information Theory
- Conditional Mutual Information and more**

Patrick Forré

Joint Entropy & Conditional Entropy

- Let X, Y be two random variables with joint distribution $p(x, y)$.
- The **Joint Entropy** of X, Y is then:

$$H(X, Y) := H((X, Y)) = H(p(x, y)) = E_p[-\log p(x, y)]$$

- The **Conditional Entropy** of X given Y is:

$$\begin{aligned} H(X|Y) &:= E_y[H(p(x|y=y))] &= \int p(y) \left[\int p(x|y) \cdot (-\log p(x|y)) dx \right] dy \\ &= H(X, Y) - H(Y) &= \int p(x,y) (-\log p(x|y)) dx dy \\ &= H(Y|X) + H(X) - H(Y) \end{aligned}$$

(Conditional) Mutual Information

- Let X, Y, Z be three random variables with joint distribution $p(x, y, z)$.
- The **Mutual Information** between X and Y is:

$$\begin{aligned} I(X; Y) &:= \text{KL}(p(x, y) \parallel p(x) \cdot p(y)) \\ &= H(X) - H(X|Y) \\ &= H(Y) - H(Y|X) \end{aligned}$$

- The **Conditional Mutual Information** between X and Y given Z is:

$$\begin{aligned} I(X; Y|Z) &:= \mathbb{E}_z [\text{KL}(p(x, y|z=z) \parallel p(x|z=z) \cdot p(y|z=z))] \\ &= H(X|Z) - H(X|Y, Z) \\ &= H(Y|Z) - H(Y|X, Z) \end{aligned}$$

Shannon's Theorem

- Let X, Y, Z be three random variables with joint distribution $p(x, y, z)$.

We always have the **inequality**:

$$I(X; Y|Z) \geq 0$$

- Furthermore, **equality** holds if and only if:

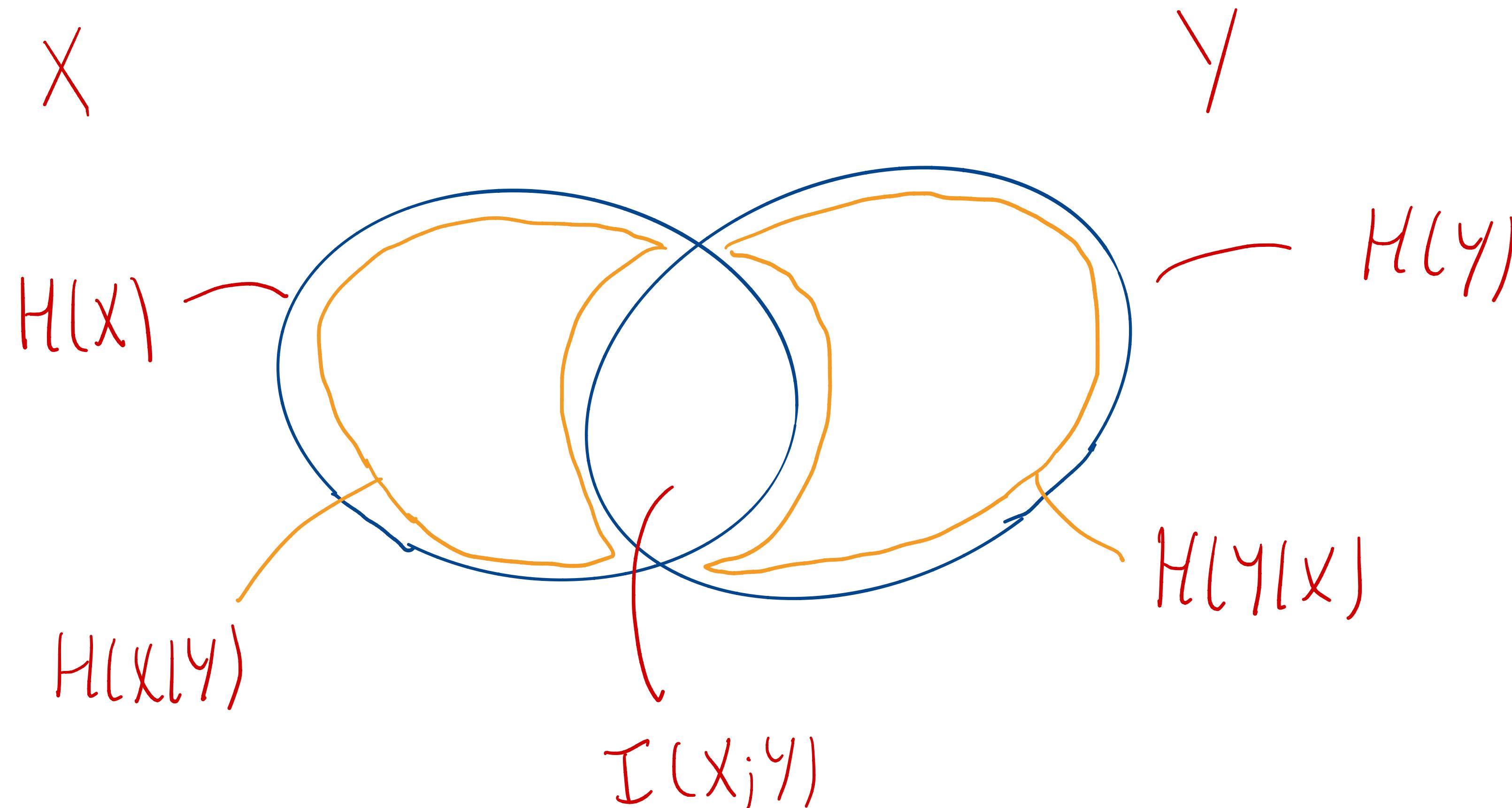
$$X \perp\!\!\!\perp Y|Z$$

- The statement also holds for the (unconditional) mutual information and (unconditional) independence (i.e. with $Z = \emptyset$).

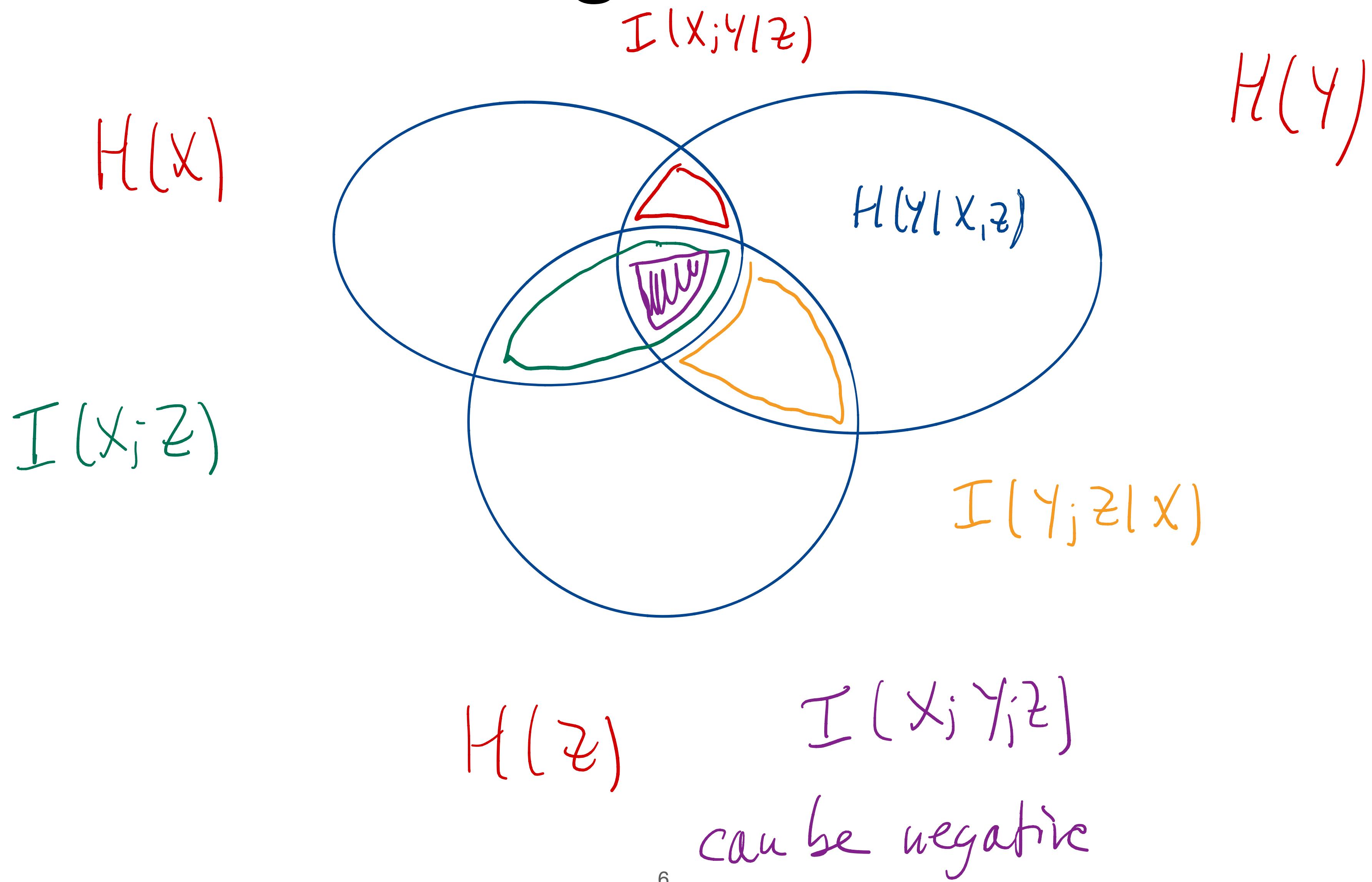
- Proof: $I(X; Y|Z) = E_Z [KL(p(X,Y|z) || p(X|z)p(Y|z))] \geq 0$

$$\stackrel{def}{=} \Leftrightarrow \forall z : KL(p(X,Y|z) || p(X|z)p(Y|z)) = 0 \Leftrightarrow p(X,Y|z) = p(X|z) \cdot p(Y|z)$$

Information Diagram for 2 Variables



Information Diagram for 3 Variables



Chain Rule

X, Y, Z, W

$$I(X; (Y, W) | Z) = I(X; Y | Z) + I(X; W | Y, Z).$$

VI
O

VI
O

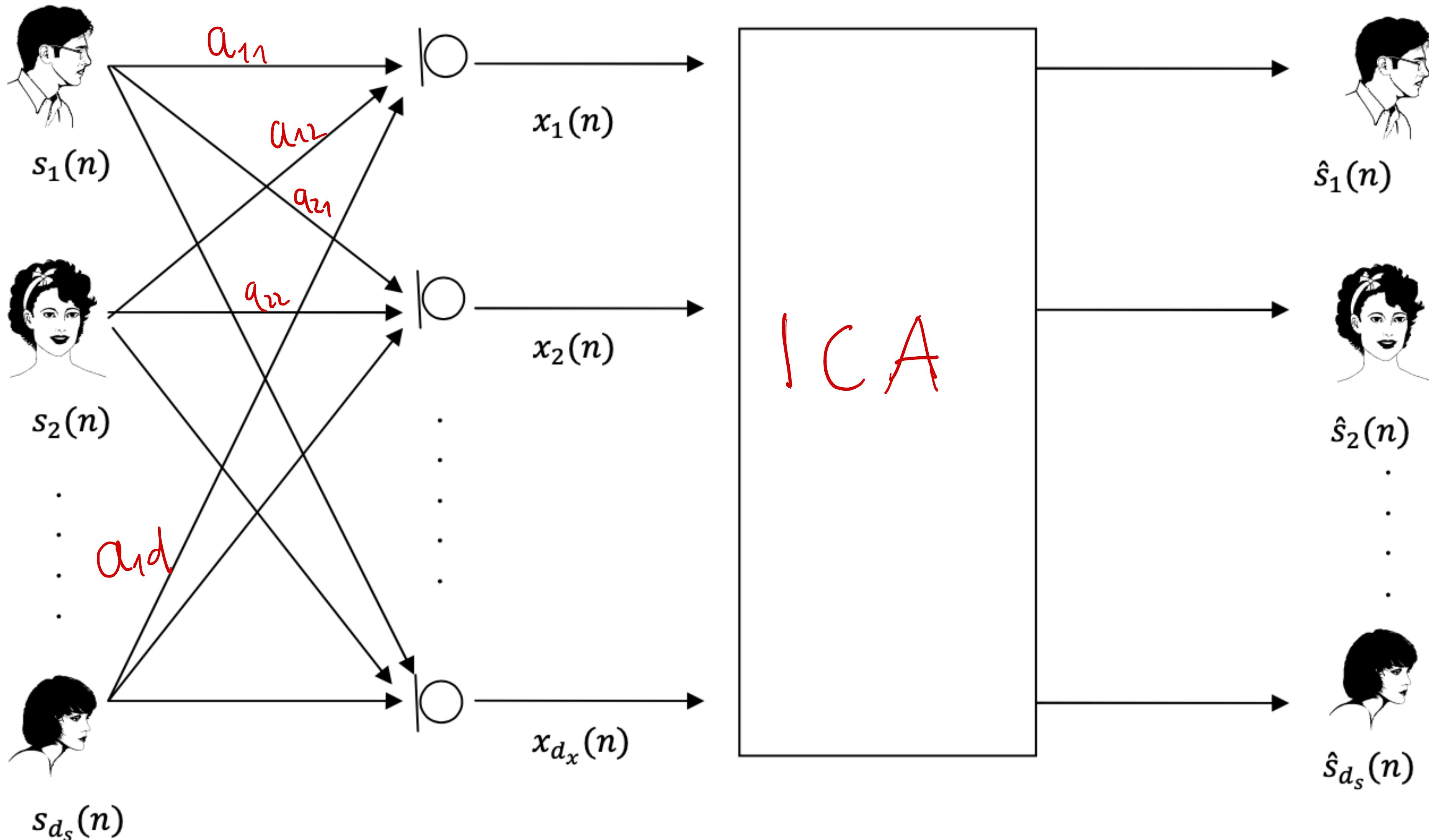
VI
O

Machine Learning 2

Independent Component Analysis (ICA)
- Objective and Assumptions

Patrick Forré

The Cocktail Party Problem



Source: S.R. Lakmal: Blind source separation in real time

ICA - Assumptions

- **Independent Components:** We have K random sources:
 $S(t) := [s_1(t), \dots, s_K(t)]^\top$ over time steps $t = 1, \dots, T$; ($K \times T$)-matrix S ;
satisfying: $s_k(t) \perp\!\!\!\perp s_{-k}(t)$ for all $k = 1, \dots, K, t = 1, \dots, T$.
- **Linear and noiseless measurements:** We measure M signals:
 $X(t) := [x_1(t), \dots, x_M(t)]^\top$; ($M \times T$)-matrix X ; such that:
$$X = A S$$
with **time-independent** ($M \times K$)-matrix A .
- **Completeness:** A is **invertible**, $M = K$, with invers $W = A^{-1}$.
- **Non-Gaussianity:** None (or at most one) of the sources $s_k(t)$ for $k = 1, \dots, K$ is Gaussian distributed, e.g. if all sources are super-Gaussian distributed.

ICA - Identifiability Theorem

- **Problem:** Recover both mixing matrix A and sources S from the data X .
- **Identifiability Theorem:** Under the mentioned assumptions, one can recover the sources $s_k(t)$ up to sign, scale and permutation and corresponding A .
- W.l.o.g. we fix the **scale** by assuming $\text{var}(s_k(t)) = 1$ for all k, t .
- Together with the independence we then have: $\text{cov}(s_i(t), s_j(t)) = \delta_{i,j}$

ICA - helpful pre-processing steps

- In most ICA algorithms the following pre-processing steps of the data are helpful (but for the algorithm from this lecture not necessary):

- **Center the data** by subtracting mean:

$$X = A S \text{ implies: } \mathbb{E}[X] = A \mathbb{E}[S] \text{ implies: } X - \mathbb{E}[X] = A(S - \mathbb{E}[S]).$$

So w.l.o.g.: $\mathbb{E}[S] = 0$ and $\mathbb{E}[X] = 0$.

- **“Whiten” the data** by using full-rank PCA and rescaling of the PCs:

So we get matrix B with BB^\top diagonal s.t. BX has covariance matrix I .

So we write X for BX and A for BA again. Then we have:

$$I = \mathbb{E}[XX^\top] = A\mathbb{E}[SS^\top]A^\top = AA^\top, \text{ thus } A \text{ is orthogonal.}$$

So we reduced parameters from K^2 to $K(K - 1)/2$. “PCA is half-ICA”.

Machine Learning 2

Independent Component Analysis (ICA)

**- Maximum Likelihood Derivation of ICA
with Natural Gradients**

Patrick Forré

ICA - Maximum Likelihood Objective

- Under the mentioned ICA assumptions: $X = AS$ or equiv.: $S = WX$
- By independence we have: $p_S(s_1, \dots, s_K) = p_{S_1}(s_1) \cdots p_{S_K}(s_K)$
- Parameters are the unknowns: $\theta = (W, p_{S_1}, \dots, p_{S_K})$
- Statistical model for x : $p(x | \theta) = p_S(W \cdot x) \cdot |\frac{\partial S}{\partial x}| = p_S(W \cdot x) \cdot |W|$
 $= p_{S_1}(w_1^\top x) \cdots p_{S_K}(w_K^\top x) \cdot |W|$
- Log-likelihood: $L(\theta) = \hat{\mathbb{E}}_x[\log p(x | \theta)] = \hat{\mathbb{E}} \left[\sum_{k=1}^K \log p_{S_k}(w_k^\top x) \right] + \log |W|$
- MLE objective: $\arg \max_{\theta} L(\theta)$

Recap: Natural Gradient

- Objective: Minimize $L(w)$
- Gradient Descent: $w^{\text{new}} := w - \alpha \cdot \nabla_w L(w),$
- Newton-Raphson: $w^{\text{new}} := w - \alpha \cdot (\nabla_w^2 L(w))^{-1} \cdot \nabla_w L(w),$
- Natural Gradient: $w^{\text{new}} := w - \alpha \cdot M \cdot \nabla_w L(w),$

where M is a so called pre-conditioner, usually a symmetric, positive-definite matrix reflecting the geometry of the space.

ICA - Gradient of Log-Likelihood

- Maximize: $L(\theta) = \hat{\mathbb{E}}_x \left[\sum_{k=1}^K \log p_{S_k}(w_k^\top x) + \log |W| \right]$

- $\frac{\partial L}{\partial w_{ij}} = \hat{\mathbb{E}}_x \left[\sum_{n=1}^N \frac{\partial}{\partial w_{ij}} \log p_{S_n}(w_n^\top x) \right] + \frac{\partial}{\partial w_{ij}} \log |W|$

$$= \hat{\mathbb{E}}_x \left[\sum_{n=1}^N \phi_n(w_n^\top x) \cdot (\delta_{ni} \cdot x_j) \right] + (W^{-1})_{ji}$$

$$\phi_n(s) := \frac{\partial}{\partial s} \log p_{S_n}(s)$$

$$= \hat{\mathbb{E}}_x [\phi_i(w_i^\top x) \cdot x_j] + (W^{-1})_{ji}$$

ICA - Activation Functions Approximation

- Approximation of the quantities related to p_{S_k} :

$$\bullet \phi_k(s) = \frac{\partial}{\partial s} \log p_{S_k}(s) \approx \begin{cases} -\tanh(s) & , s_k \text{ super-Gaussian} \\ \tanh(s) - 5 & , s_k \text{ sub-Gaussian} \end{cases}$$

ICA - Natural Gradient

- A good data-independent preconditioner for $L(\theta)$ turns out to be:

$$M_{(lm),(ij)} = \sum_v W_{vm} W_{vj} \cdot \delta_{ij}$$

Natural gradient: $(M \cdot \nabla_W L)_{(lm)} = \sum_{i,j,v} W_{vm} W_{vj} \delta_{ij} (\nabla_W L)_{(ij)}$

$$= \sum_{j,v} W_{vm} W_{vj} \cdot [\hat{E}_x [\phi_e(w_e^T x) \cdot x_j] + (W^{-1})_{je}]$$

$$= \sum_v W_{vm} \hat{E}_x [\phi_e(w_e^T x) (w_v^T \cdot x)] + \sum_v W_{vm} \delta_{ve}$$

$$= w_{em} + \sum_v W_{vm} \hat{E}_x [\phi_e(w_e^T x) \cdot (w_v^T \cdot x)].$$

ICA - Update Rule

$x(t)$

- $W^{\text{new}} = W + \alpha \cdot [W + \phi(W \cdot x(t)) \cdot x(t)^T W^T W]$

$$\hat{W} \sim MLE \quad (\text{after convergence})$$

$$\hat{A} \simeq \hat{W}^{-1}$$

- After convergence reconstruct:

$$\hat{s}_k(t) := (\hat{W} \cdot x(t))_k$$

$$\hat{s}(t) := \hat{W} \cdot x(t)$$

Machine Learning 2

Independent Component Analysis (ICA)

- Covariant Online ICA Algorithm**
- Summary**

Patrick Forré

ICA - Assumptions for Identifiability

- **Independent sources:** $s_k(t) \perp\!\!\!\perp s_{-k}(t)$.
- **Linear and noiseless measurements:** $X = A S$ with matrix A .
- **Completeness:** A is invertible, $W := A^{-1}$.
- **Non-Gaussianity:** Sources $s_k(t)$ are non-Gaussian (or at most one is).
- **Theorem:** Under those assumptions both mixing matrix A and sources $s_k(t)$ can be identified up to sign, scale and permutation of the sources.

Covariant Online ICA Algorithm

- Choose learning rate α , activation functions ϕ_k (usually $\phi_k(s) = -\tanh(s)$ for super-Gaussian and $\phi_k(s) = \tanh(s) - s$ for sub-Gaussian), initialize W
- Until convergence take new data point $x(t)$ and do:
 - $c(t) := W \cdot x(t)$
 - $z(t) := \phi(c(t))$
 - $y(t) := W^\top \cdot c(t)$
 - $W^{\text{new}} := W + \alpha \cdot [W + z(t) \cdot y(t)^\top]$
 - Reconstruct sources: $\hat{s}(t) := \hat{W} \cdot x(t)$