

CS340: MACHINE LEARNING

REVIEW FOR MIDTERM

KEVIN MURPHY

- We want to make models of data so we can find patterns and predict the future.
- We will use probability theory to represent our uncertainty about what the right model is; this will induce uncertainty in our predictions.
- We will update our beliefs about the model as we acquire data (this is called “learning”); as we get more data, the uncertainty in our predictions will decrease.

---

OUTLINE

- Probability basics
- Probability density functions
- Parameter estimation
- Prediction
- Multivariate probability models
- Conditional probability models
- Information theory

---

PROBABILITY BASICS

- Chain (product) rule

$$p(A, B) = p(A)p(B|A) \tag{1}$$

- Marginalization (sum) rule

$$p(B) = \sum_a p(A = a, B) \tag{2}$$

- Hence we derive Bayes rule

$$p(A|B) = \frac{p(A, B)}{p(B)} \tag{3}$$

## PMF

---

- A probability mass function (PMF) is defined on a discrete random variable  $X \in \{1, \dots, K\}$  and satisfies

$$1 = \sum_{x=1}^K P(X = x) \quad (4)$$

$$0 \leq P(X = x) \leq 1 \quad (5)$$

- This is just a histogram!

## MOMENTS OF A DISTRIBUTION

---

- Mean (first central moment)

$$\mu = E[X] = \sum_x xp(x) \quad (8)$$

- Variance (second central moment)

$$\sigma^2 = \text{Var} [X] = \sum_x (x - \mu)^2 p(x) = E[X^2] - \mu^2 \quad (9)$$

## PDF

---

- A probability density function (PDF) is defined on a continuous random variable  $X \in \mathbb{R}$  or  $X \in \mathbb{R}^+$  or  $X \in [0, 1]$  etc and satisfies

$$1 = \int_{\mathcal{X}} p(X = x) dx \quad (6)$$

$$0 \leq p(X = x) \quad (7)$$

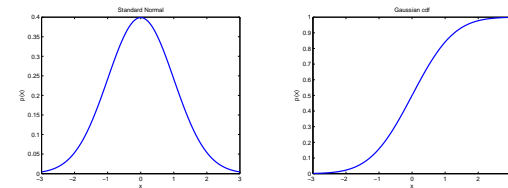
- Note that it is possible that  $p(X = x) > 1$ , since multiplied by  $dx$ .
- $P(x \leq X \leq x + dx) \approx p(x)dx$
- Statistics books often write  $f_X(x)$  for pdfs and use  $P()$  for pmf's.

## CDFs

---

The cumulative distribution function is defined as

$$F(x) = P(X \leq x) = \int_{-\infty}^x f(x) dx \quad (10)$$



## QUANTILES OF A DISTRIBUTION

---

$$P(a \leq X \leq b) = F_X(b) - F_X(a) \quad (11)$$

eg if  $Z \sim \mathcal{N}(0, 1)$ ,  $F_Z(z) = \Phi(z)$ . Then

$$p(-1.96 \leq Z \leq 1.96) = \Phi(1.96) - \Phi(-1.96) \quad (12)$$

$$= 0.975 - 0.025 = 0.95 \quad (13)$$

Hence  $X \sim \mathcal{N}(\mu, \sigma)$ ,

$$p(-1.96\sigma < X - \mu < 1.96\sigma) = 1 - 2 \times 0.025 = 0.95 \quad (14)$$

Hence the interval  $\mu \pm 1.96\sigma$  contains 0.95 mass.

## BERNOULLI DISTRIBUTION

---

- $X \in \{0, 1\}$ ,  $p(X = 1) = \theta$ ,  $p(X = 0) = 1 - \theta$
- $X \sim Be(\theta)$
- $E[X] = \theta$

## OUTLINE

---

- Probability basics ✓
- Probability density functions
- Parameter estimation
- Prediction
- Multivariate probability models
- Conditional probability models
- Information theory

## MULTINOMIAL DISTRIBUTION

---

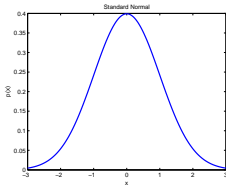
- $X \in \{1, \dots, K\}$ ,  $p(X = k) = \theta_k$ ,  $X \sim Mu(\theta)$
- $0 \leq \theta_k \leq 1$
- $\sum_{k=1}^K \theta_k = 1$
- $E[X_i] = \theta_i$

## (UNIVARIATE) GAUSSIAN DISTRIBUTION

- $X \in \mathbb{R}$ ,  $X \sim \mathcal{N}(\mu, \sigma)$ ,  $\mu \in \mathbb{R}$ ,  $\sigma \in \mathbb{R}^+$

$$\mathcal{N}(x|\mu, \sigma) \stackrel{\text{def}}{=} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} \quad (15)$$

- $E[X] = \mu$

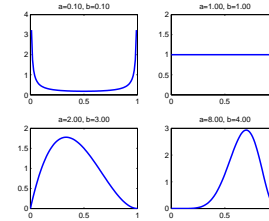


## BETA DISTRIBUTION

- $X \in [0, 1]$ ,  $X \sim Be(\alpha_1, \alpha_0)$ ,  $0 \leq \alpha_0, \alpha_1$

$$E[X] = \frac{\alpha_1}{\alpha_1 + \alpha_0}$$

$$Be(X|\alpha_1, \alpha_0) \propto [X^{\alpha_1-1}(1-X)^{\alpha_0-1}]$$

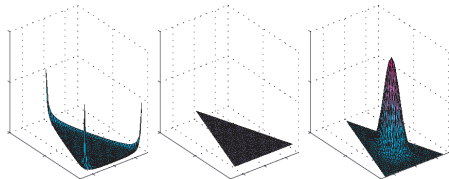
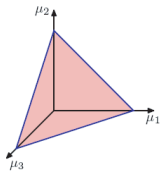


## DIRICHLET DISTRIBUTION

- $X \in [0, 1]^K$ ,  $X \sim Dir(\alpha_1, \dots, \alpha_K)$ ,  $0 \leq \alpha_k$

- $E[X_i] = \frac{\alpha_i}{\alpha}$ , where  $\alpha = \sum_k \alpha_k$ .

$$Dir(X|\alpha_1, \dots, \alpha_K) \propto \prod_k X_k^{\alpha_k-1}$$



## OUTLINE

- Probability basics ✓
- Probability density functions ✓
- Parameter estimation
- Prediction
- Multivariate probability models
- Conditional probability models
- Information theory

---

PARAMETER ESTIMATION: MLE

- Likelihood of iid data  $D = \{x_n\}_{n=1}^N$

$$L(\theta) = p(D|\theta) = \prod_n p(x_n|\theta) \quad (16)$$

- Log likelihood

$$\ell(\theta) = \log p(D|\theta) = \sum_n \log p(x_n|\theta) \quad (17)$$

- MLE = Maximum likelihood estimate

$$\hat{\theta} = \arg \max_{\theta} L(\theta) = \arg \max_{\theta} \ell(\theta) \quad (18)$$

hence

$$\frac{\partial L}{\partial \theta_i} = 0 \quad (19)$$

---

MLE FOR MULTINOMIALS

$$p(X_n|\theta) = \prod_k \theta_k^{I(X_n=k)} \quad (26)$$

$$\ell(\theta) = \sum_n \sum_k I(x_n = k) \log \theta_k \quad (27)$$

$$= \sum_k N_k \log \theta_k \quad (28)$$

$$\tilde{l} = \sum_k N_k \log \theta_k + \lambda \left( 1 - \sum_k \theta_k \right) \quad (29)$$

$$\frac{\partial \tilde{l}}{\partial \theta_k} = \frac{N_k}{\theta_k} - \lambda = 0 \quad (30)$$

$$\hat{\theta}_k = \frac{N_k}{N} \quad (31)$$

---

MLE FOR BERNOULLIS

$$p(X_n|\theta) = \theta^{X_n}(1 - \theta)^{1-X_n} \quad (20)$$

$$\ell(\theta) = \sum_n X_n \log \theta + \sum_n (1 - X_n) \log(1 - \theta) \quad (21)$$

$$= N_1 \log \theta + N_0 \log(1 - \theta) \quad (22)$$

$$N_1 = \sum_n I(X_n = 1) \quad (23)$$

$$\frac{dL}{d\theta} = \frac{N_1}{\theta} - \frac{N_0}{1 - \theta} = 0 \quad (24)$$

$$\hat{\theta} = \frac{N_1}{N_1 + N_0} \quad (25)$$

---

PARAMETER ESTIMATION: BAYESIAN

$$p(\theta|D) \propto p(D|\theta)p(\theta) \quad (32)$$

$$= \left[ \prod_n p(x_n|\theta) \right] p(\theta) \quad (33)$$

We say  $p(\theta)$  is conjugate to  $p(D|\theta)$  if  $p(\theta|D)$  in same family as  $p(\theta)$ . This allows recursive (sequential) updating.

MLE is a point estimate. Bayes computes full posterior, natural model of uncertainty (e.g., posterior variance, entropy, credible intervals, etc)

---

## PARAMETER ESTIMATION: MAP

---

For many models, as  $|D| \rightarrow \infty$ ,

$$p(\theta|D) \rightarrow \delta(\theta - \hat{\theta}_{MAP}) \rightarrow \delta(\theta - \hat{\theta}_{MLE}) \quad (34)$$

where

$$\hat{\theta}_{MAP} = \arg \max_{\theta} p(D|\theta)p(\theta) \quad (35)$$

MAP (maximum a posteriori) estimate is a posterior mode. Also called a penalized likelihood estimate since

$$\hat{\theta}_{MAP} = \arg \max_{\theta} \log p(D|\theta) - \lambda c(\theta) \quad (36)$$

where  $c(\theta) = -\log p(\theta)$  is a penalty or regularizer, and  $\lambda$  is the strength of the regularizer.

---

## BAYES ESTIMATE FOR BERNOULLI

---

$$X_n \sim Be(\theta) \quad (37)$$

$$\theta \sim Beta(\alpha_1, \alpha_0) \quad (38)$$

$$p(\theta|D) \propto \left[ \prod_n \theta^{X_n} (1-\theta)^{1-X_n} \right] \theta^{\alpha_1-1} (1-\theta)^{\alpha_0-1} \quad (39)$$

$$= \theta^{N_1} (1-\theta)^{N_0} \theta^{\alpha_1-1} (1-\theta)^{\alpha_0-1} \quad (40)$$

$$\propto Beta(\alpha_1 + N_1, \alpha_0 + N_0) \quad (41)$$

---

## BAYES ESTIMATE FOR MULTINOMIAL

---

$$X_n \sim Mu(\theta) \quad (42)$$

$$\theta \sim Dir(\alpha_1, \dots, \alpha_K) \quad (43)$$

$$p(\theta|D) \propto \left[ \prod_n \prod_k \theta_k^{I(X_n=k)} \right] \prod_k \theta_k^{\alpha_k-1} \quad (44)$$

$$= \prod_k \theta_k^{N_k} \theta_k^{\alpha_k-1} \quad (45)$$

$$\propto Dir(\alpha_1 + N_1, \dots, \alpha_K + N_K) \quad (46)$$

---

## OUTLINE

---

- Probability basics ✓
- Probability density functions ✓
- Parameter estimation ✓
- Prediction
- Multivariate probability models
- Conditional probability models
- Information theory

## PREDICTING THE FUTURE

---

Posterior predictive density gotten by Bayesian model averaging

$$p(X|D) = \int p(X|\theta)p(\theta|D)d\theta \quad (47)$$

Plug-in principle: if  $p(\theta|D) \approx \delta(\theta - \hat{\theta})$ , then

$$p(X|D) \approx p(X|\hat{\theta}) \quad (48)$$

## CROSS VALIDATION

---

We can use CV to find the model with best predictive performance.

## PREDICTIVE FOR BERNOULLIS

---

$$p(\theta|D) = \text{Beta}(\alpha_1 + N_1, \alpha_0 + N_0) \quad (49)$$

$$\stackrel{\text{def}}{=} \text{Beta}(\alpha'_1, \alpha'_0) \quad (50)$$

$$p(X = 1|D) = \int p(X = 1|\theta)p(\theta|D)d\theta \quad (51)$$

$$= \int \theta p(\theta|D)d\theta \quad (52)$$

$$= E[\theta|D] \quad (53)$$

$$= \frac{\alpha'_1}{\alpha'_1 + \alpha'_0} \quad (54)$$

## OUTLINE

---

- Probability basics ✓
- Probability density functions ✓
- Parameter estimation ✓
- Prediction ✓
- Multivariate probability models
- Conditional probability models
- Information theory

---

MULTIVARIATE PROBABILITY MODELS

- Each data case  $X_n$  is now a vector of variables  $X_{ni}$ ,  $i = 1 : p$ ,  $n = 1 : N$
- e.g.,  $X_{ni} \in \{1, \dots, K\}$  is the  $i$ 'th word in  $n$ 'th sentence
- e.g.,  $X_{ni} \in \mathbb{R}$  is the  $i$ 'th attribute in the  $n$ 'th feature vector
- We need to define  $p(X_n|\theta)$  for feature vectors of potentially variable size.

---

INDEPENDENT FEATURES (BAG OF WORDS MODEL)

$$p(X_n|\theta) = \prod_i p(X_{ni}|\theta_i) \quad (55)$$

e.g.,  $X_{ni} \in \{0, 1\}$ , product of Bernoullis

$$p(X_n|\theta) = \prod_i \theta_i^{X_{ni}} (1 - \theta_i)^{1 - X_{ni}} \quad (56)$$

e.g.,  $X_{ni} \in \{1, \dots, K\}$ , product of multinomials

$$p(X_n|\theta) = \prod_i \prod_k \theta_{ik}^{I(X_{ni}=k)} \quad (57)$$

---

MLE FOR BAG OF WORDS MODEL

$$p(D|\theta) = \prod_n \prod_i p(X_{ni}|\theta_i) \quad (58)$$

We can estimate each  $\theta_i$  separately, eg. product of multinomials

$$p(D|\theta) = \prod_n \prod_i \prod_k \theta_{ik}^{I(X_{ni}=k)} \quad (59)$$

$$= \prod_i \prod_k \theta_{ik}^{N_{ik}} \quad (60)$$

$$N_{ik} \stackrel{\text{def}}{=} \sum_n I(X_{ni} = k) \quad (61)$$

$$\ell(\theta_i) = \sum_k N_{ik} \log \theta_{ik} \quad (62)$$

$$\theta_{ik} = \frac{N_{ik}}{N_i} \quad (63)$$

---

BAYES ESTIMATE FOR BAG OF WORDS MODEL

$$p(\theta|D) \propto \prod_i p(\theta_i) [\prod_n \prod_i p(X_{ni}|\theta_i)] \quad (64)$$

We can update each  $\theta_i$  separately, eg. product of multinomials

$$p(\theta|D) = \prod_i \text{Dir}(\theta_i|\alpha_{i1}, \dots, \alpha_{iK}) [\prod_k \theta_{ik}^{N_{ik}}] \quad (65)$$

$$\propto \prod_i \text{Dir}(\theta_i|\alpha_{i1} + N_{i1}, \dots, \alpha_{iK} + N_{iK}) \quad (66)$$

Factored prior  $\times$  factored likelihood = factored posterior



$$p(X_n|\theta) = p(X_{n1}|\pi) \prod_{t=2}^n p(X_{nt}|X_{n,t-1}, A) \quad (67)$$

For a discrete state space,  $X_{nt} \in \{1, \dots, K\}$ ,

$$\pi_i = p(X_1 = i) \quad (68)$$

$$A_{ij} = p(X_t = j | X_{t-1} = i) \quad (69)$$

$$p(X_n|\theta) = \prod_i \pi_i^{I(X_{n1}=i)} \prod_t \prod_i \prod_j A_{ij}^{I(X_{nt}=j, X_{n,t-1}=i)} \quad (70)$$

$$= \prod_i \pi_i^{N_i^1} \prod_i \prod_j A_{ij}^{N_{ij}} \quad (71)$$

$$N_i^1 \stackrel{\text{def}}{=} I(X_{n1} = i) \quad (72)$$

$$N_{ij} \stackrel{\text{def}}{=} \sum_t I(X_{n,t} = j, X_{n,t-1} = i) \quad (73)$$

$$p(D|\theta) = \prod_n \prod_i \pi_i^{I(X_{n1}=i)} \prod_t \prod_i \prod_j A_{ij}^{I(X_{nt}=j, X_{n,t-1}=i)} \quad (74)$$

$$= \prod_i \pi_i^{N_i^1} \prod_i \prod_j A_{ij}^{N_{ij}} \quad (75)$$

$$N_i^1 \stackrel{\text{def}}{=} \sum_n I(X_{n1} = i) \quad (76)$$

$$N_{ij} \stackrel{\text{def}}{=} \sum_n \sum_t I(X_{n,t} = j, X_{n,t-1} = i) \quad (77)$$

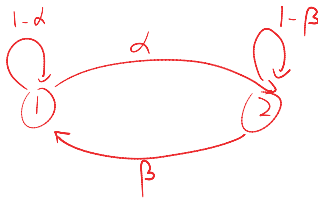
$$\hat{A}_{ij} = \frac{N_{ij}}{\sum_{j'} N_{ij'}} \quad (78)$$

$$\hat{\pi}_i = \frac{N_i^1}{N} \quad (79)$$

STATIONARY DISTRIBUTION FOR A MARKOV CHAIN

$\pi_i$  is fraction of time we spend in state  $i$ , given by principle eigenvector

$$A\pi = \pi \quad (80)$$



Balance flow across cut set

$$\pi_1 \alpha = \pi_2 \beta \quad (81)$$

Since  $\pi_1 + \pi_2 = 1$ , we have

$$\pi_1 = \frac{\beta}{\alpha + \beta}, \quad \pi_2 = \frac{\alpha}{\alpha + \beta}, \quad (82)$$

OUTLINE

- Probability basics ✓
- Probability density functions ✓
- Parameter estimation ✓
- Prediction ✓
- Multivariate probability models ✓
- Conditional probability models
- Information theory

## CONDITIONAL DENSITY ESTIMATION

---

• Goal: learn  $p(y|x)$  where  $x$  is input and  $y$  is output.

• Generative model

$$p(y|x) \propto p(x|y)p(y) \quad (83)$$

• Discriminative model e.g., logistic regression

$$p(y|x) = \sigma(w^T x) \quad (84)$$

## NAIVE BAYES = CONDITIONAL BAG OF WORDS MODEL

---

Generative model

$$p(y|x) \propto p(x|y)p(y) \quad (85)$$

in which features of  $x$  are conditionally independent given  $y$ :

$$p(y|x) \propto p(y) \left[ \prod_{i=1}^p p(x_i|y) \right] \quad (86)$$

e.g., for multinomials

$$p(y|x) \propto \prod_c \pi_c^{I(y=c)} \left[ \prod_i \prod_k \theta_{cik}^{I(x_i=k, y=c)} \right] \quad (87)$$

## MLE FOR NAIVE BAYES

---

We estimate params for each feature separately, partitioning the data based on label  $y$ .

$$p(D|\theta, \pi) = \left[ \prod_c \pi_c^{N_c} \right] \prod_i \left[ \prod_k \prod_c \theta_{cik}^{N_{cik}} \right] \quad (88)$$

$$N_{cik} = \sum_n I(x_{ni} = k, y_n = c) \quad (89)$$

$$\hat{\theta}_{cik} = \frac{N_{cik}}{\sum_{k'} N_{cik'}} \quad (90)$$

$$\hat{\pi}_c = \frac{N_c}{N} \quad (91)$$

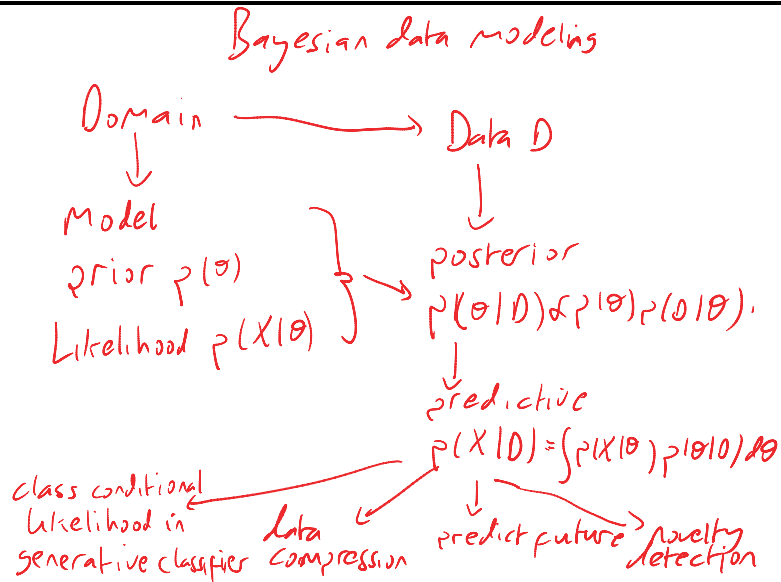
## CONDITIONAL MARKOV MODEL

---

$$p(x|y) = p(x_1|y) \prod_t p(x_t|x_{t-1}, y) \quad (92)$$

Fit a separate Markov model for every class of  $y$ . See homework 5.

SUMMARY



SUMMARY

- Probability basics: Bayes rule, pdfs, pmfs
- Probability density functions: Bernoulli, Multinomial; Gaussian, Beta, Dirichlet
- Parameter estimation: MLE, Bayesian
- Prediction: Plug-in, Bayesian, cross validation
- Multivariate probability models: Bag of words, Markov models
- Conditional probability models: Conditional BOW (Naive Bayes), Conditional Markov
- Information theory

ENTROPY

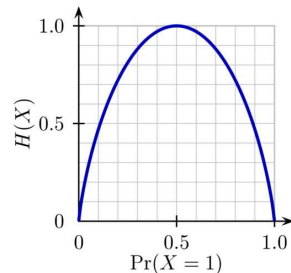
- Measure of uncertainty.
- Definition for PMF  $p_k, k = 1 : K$

$$H(p) = - \sum_k p_k \log_2 p_k \quad (93)$$

- eg binary entropy function,  $p_1 = \theta, p_0 = 1 - \theta$

$$H(\theta) = -[p(X = 1) \log_2 p(X = 1) + p(X = 0) \log_2 p(X = 0)] \quad (94)$$

$$= -[\theta \log_2 \theta + (1 - \theta) \log_2 (1 - \theta)] \quad (95)$$



JOINT ENTROPY

$$H(X, Y) = - \sum_{x,y} p(x, y) \log p(x, y) \quad (96)$$

If  $X \perp Y$ , then  $H(X, Y) = H(X) + H(Y)$ .

Pf:

$$H(X, Y) = - \sum_{x,y} p(x)p(y) \log p(x)p(y) \quad (97)$$

$$= - \sum_{x,y} p(x)p(y) \log p(x) - \sum_{x,y} p(x)p(y) \log p(y) \quad (98)$$

In general

$$H(X, Y) \leq H(X) + H(Y) \quad (99)$$

Pf:

$$I(X, Y) = H(X) + H(Y) - H(X, Y) \geq 0 \quad (100)$$

CONDITIONAL ENTROPY

$$H(Y|X) \stackrel{\text{def}}{=} \sum_x p(x)H(Y|X = x) \quad (101)$$

$$= H(X, Y) - H(X) \quad (102)$$

In hw 5, you showed that  $H(Y|X) \leq H(Y)$ .

Pf

$$H(Y|X) = H(X, Y) - H(X) \leq H(X) + H(Y) - H(X) \quad (103)$$

But  $\exists y. H(X|y) > H(X)$ .

MUTUAL INFORMATION

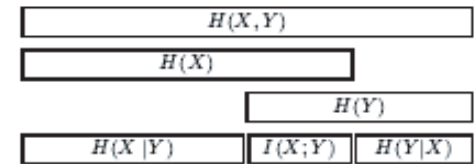
$$I(X; Y) = \sum_{y \in Y} \sum_{x \in X} p(x, y) \log \frac{p(x, y)}{p(x)p(y)} \quad (104)$$

In hw 5, you showed that

$$I(X, Y) = H(X) - H(X|Y) = H(Y) - H(Y|X) \quad (105)$$

Substituting  $H(Y|X) = H(X, Y) - H(X)$  yields

$$I(X, Y) = H(X) + H(Y) - H(X, Y) \quad (106)$$



MUTUAL INFORMATION  $\geq 0$

$I(X, Y) = 0$  if  $X \perp Y$ . Pf.

$$I(X; Y) = \sum_{y \in Y} \sum_{x \in X} p(x)p(y) \log \frac{p(x)p(y)}{p(x)p(y)} \quad (107)$$

$$= \sum_y \sum_x p(x)p(y) \log 1 = 0 \quad (108)$$

$I(X, Y) \geq 0$ . Pf.

$$I(X, Y) = KL(p(x, y)||p(x)p(y)) \geq 0 \quad (109)$$

RELATIVE ENTROPY

$$\mathcal{D}(p||q) \stackrel{\text{def}}{=} \sum_k p_k \log \frac{p_k}{q_k} \quad (110)$$

$$= \sum_k p_k \log p_k - \sum_k p_k \log q_k \quad (111)$$

$$= - \sum_k p_k \log q_k - H(p) \quad (112)$$

Hence measures extra number of bits needed to encode data if we use  $q$  instead of  $p$ .

$\mathcal{D}(p||q) = 0$  if  $p = q$ . Pf: trivial.

$\mathcal{D}(p||q) \geq 0$ . Pf: use Jensen's inequality.

## MDL PRINCIPLE

---

Minimum description length principle.

To encode an event  $X = x$  which occurs with  $p(x)$  takes  $-\log p(x)$  bits. Total cost of data  $D$  and model  $H$  is

$$L(D, H) = -\log p(D|H) - \log p(H) \quad (113)$$

MDL principle says: pick

$$H^* = \arg \min_H L(D, H) \quad (114)$$

This is equivalent to the MAP hypothesis

$$H^* = \arg \max_H \log p(D|H) + \log p(H) \quad (115)$$

