

# PROBABILITY and UNCERTAINTY

ALEATORIC UNCERTAINTY

EPISTEMIC UNCERTAINTY

- FREQUENTIST INTERPRETATION

$P(X=i)$  probability of observing outcome  $i$   
of an experiment

- BAYESIAN INTERPRETATION

$P(X=i)$  belief of observing  $i$

# DISCRETE PROBABILITY DISTRIBUTIONS

$$P(X=i) \in [0,1] \quad \mathcal{X} \text{ SET OF OUTCOMES}$$

$$\sum_{i \in \mathcal{X}} P(X=i) = 1$$

- RANDOM VARIABLE  $X \in \mathcal{X}$ ,  $P(X=i)$

$$P(X)$$

$X, Y$  RANDOM VARIABLES

$$P(X, Y) \quad \text{joint p.d.} \quad P(X=i, Y=j)$$

$$P(X|Y) \quad \text{conditional probability of } X \text{ given } Y \text{ ( } Y \text{ has a fixed value)}$$

$$P(X) \quad \text{marginal p.d. of } X$$

BERNOULLI  $x \in \{0, 1\}$   $p(x) = \theta \in [0, 1]$   
COIN  $\theta = \frac{1}{2}$

BINOMIAL DISTRIBUTION  $n$  THROWS

$x \in \{0, \dots, n\}$   $p(x=i) = \binom{n}{i} \theta^i (1-\theta)^{n-i}$

# LAWS OF PROBABILITY

• NORMALIZATION  $\sum_{x \in \mathcal{X}} P(x) = 1$   $\left( \sum_{i \in \mathcal{X}} P(x=i) = 1 \right)$

• SUM RULE  $P(x) = \sum_{y \in \mathcal{Y}} P(x, y)$

• PRODUCT RULE  $P(x, y) = P(x|y)P(y)$

$$P(x|y) = \frac{P(x, y)}{P(y)}$$

• BAYES THEOREM

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)}$$

$$P(x) = \sum_{y \in \mathcal{Y}} P(x|y)P(y)$$

## INDEPENDENCE

X and y are independent

$$\begin{aligned} P(x, y) &= P(x) P(y) \\ \parallel & \\ P(x|y) P(y) & \Rightarrow P(x|y) = P(x) \end{aligned}$$

# CONTINUOUS PROBABILITY DISTRIBUTIONS

$$\mathcal{X} \subseteq \mathbb{R}^n \cdot n \geq 1$$

WE CONSIDER DISTRIBUTIONS HAVING A DENSITY (in  $\mathbb{R}$ )

$$p: \mathcal{X} \rightarrow \mathbb{R} \text{ s.t. } p(x) \geq 0 \text{ and } \int_{\mathcal{X}} p(x) dx = 1$$

$$p(x \in A) = p(A) = \int_A p(x) dx$$

A is a nice subset (e.g.  $A = [a, b] \subseteq \mathbb{R}$ )

- $\int_{\mathcal{X}} p(x) dx = 1$
- $p(x) = \int_{\mathcal{Y}} p(x, y) dy$
- $p(x, y) = \underbrace{p(x|y)} p(y)$

# EXPECTATIONS

$$\mathcal{X} \ni x \quad f: \mathcal{X} \rightarrow \mathbb{R}$$

$$E[f] = E[f(x)] = \int_{\mathcal{X}} f(x) p(x) dx \quad \left( \begin{array}{l} \text{discrete} \\ \downarrow \\ = \sum_{x \in \mathcal{X}} f(x) \cdot p(x) \end{array} \right)$$

$$x \in \{0, 1\} \quad f(0) = -1000, \quad f(1) = +100, \quad p(x=1) = 0.75$$

$$E[f(x)] = -1000 \cdot (0.25) + 100 \cdot (0.75) = -175$$

$$\mu = E[x] \text{ or mean}$$

$$\text{VAR} = E[(x - \mu)^2] \text{ or variance}$$

$$\text{COV}(x, y) = E[(x - \mu_x)(y - \mu_y)]$$

$$\text{CORR}(x, y) = \frac{\text{COV}(x, y)}{\sqrt{\text{VAR}(x) \text{VAR}(y)}}$$

$\uparrow$   
[ -1, 1 ]

$$E[f(x)] = \int_x f(x) p(x) dx$$



$X \rightarrow$  generate SAMPLES  $x_1, \dots, x_N$  by using some sampling algorithm.

$$E[f(x)] \approx \underbrace{\frac{1}{N} \sum_{i=1}^N f(x_i)}_{\downarrow N \rightarrow \infty} \\ E[f(x)]$$



# PROBABILITIES - THE FORMAL WAY

## $\sigma$ -ALGEBRA

$X$  be a set,  $\mathcal{S} \subseteq 2^X$  is a  $\sigma$ -ALGEBRA

1)  $\emptyset, X \in \mathcal{S}$

2)  $A \in \mathcal{S} \Rightarrow A^c \in \mathcal{S}$  ( $A^c = X \setminus A$ )

3)  $A_n \in \mathcal{S} \Rightarrow \bigcup_n A_n \in \mathcal{S}, n \in \mathbb{N}$

$(X, \mathcal{S})$  is a MEASURABLE SPACE.

$X = \mathbb{R}^n$ ,  $\mathcal{B} \subseteq 2^X$  BOREL  $\sigma$ -ALGEBRA, the smallest

$\sigma$ -ALGEBRA containing open sets of the standard topology of  $\mathbb{R}^n$

$(X, \mathcal{A}), (Y, \mathcal{B})$  measurable spaces

$f: (X, \mathcal{A}) \rightarrow (Y, \mathcal{B})$  is measurable

$$f^{-1}(B) \in \mathcal{A}, B \in \mathcal{B}$$

### PROBABILITY MEASURE

$(X, \mathcal{A})$  measurable space. A p.m. on  $(X, \mathcal{A})$  is

$$p: \mathcal{A} \rightarrow [0, 1] \text{ s.t.}$$

$$p(\emptyset) = 0$$

$$p(A^c) = 1 - p(A)$$

if  $A_n \in \mathcal{A}$  disjoint  $(A_i \cap A_j) = \emptyset$

$$p\left(\bigcup_n A_n\right) = \sum_n p(A_n)$$

# RANDOM VARIABLES

$$(\Omega, \mathcal{F})$$

SAMPLE SPACE

$\omega \in \Omega$  state of the world

$$P_{\Omega}(S), S \in \mathcal{F}$$

$$(\mathcal{X}, \mathcal{A}) \text{ measurable spaces}$$

SPACE of MEASURES

$x: \Omega \rightarrow \mathcal{X}$  measurable function

$$A \in \mathcal{A} \Rightarrow x^{-1}(A) \in \mathcal{F}$$

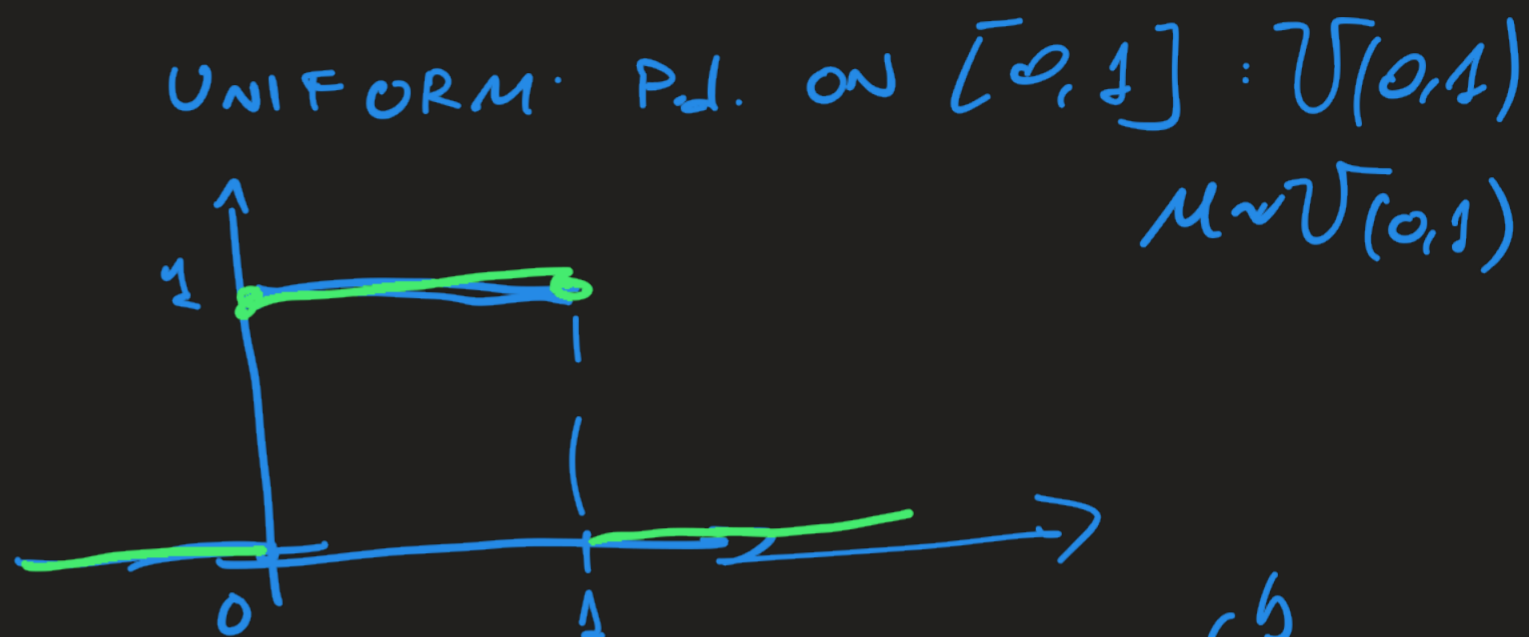
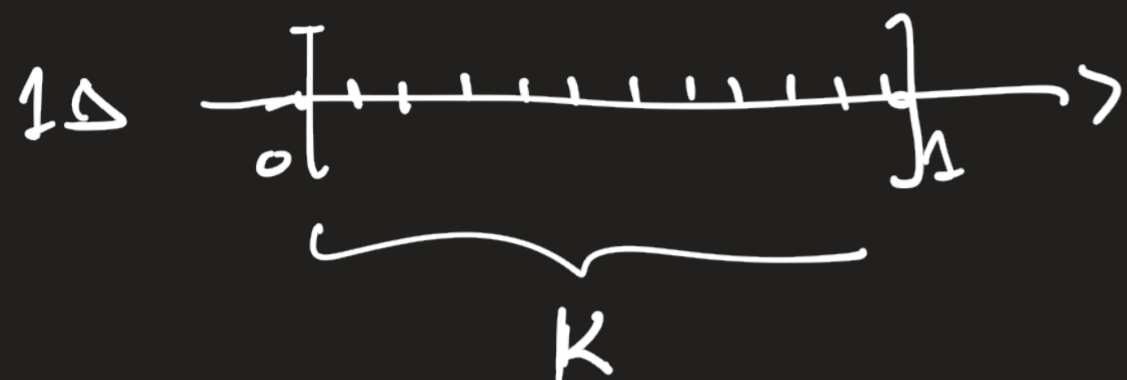
$$P_x(A) = P_{\Omega}(x^{-1}(A))$$

$$P(x \in A)$$

$$\rightarrow \int_A p(x) dx$$

density of  $x$

# CURSE OF DIMENSIONALITY



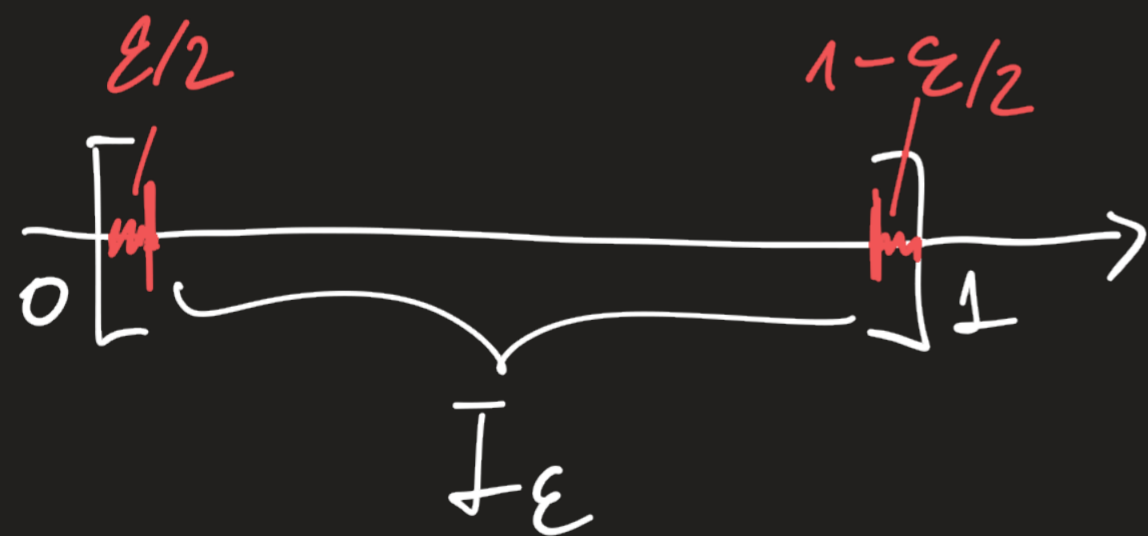
$n$ -dimensional, a grid for  $[0,1]^n$

will have  $K^n$  points

For  $n$  large, then  $K^n$  is out of reach

$$P(\mu \in [a,b]) = b-a = \int_a^b 1 d\mu$$

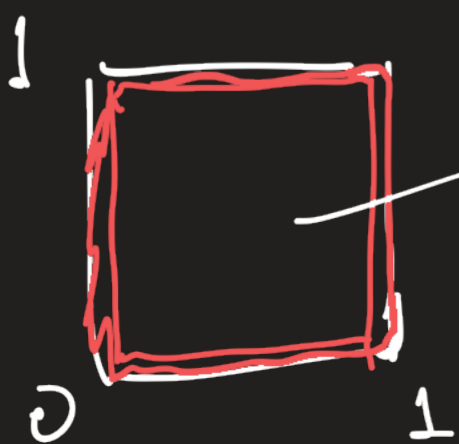
$0 \leq a \leq b \leq 1$



$$P(\mu \notin \underbrace{[\frac{\epsilon}{2}, 1-\frac{\epsilon}{2}]}_{I_\epsilon}) = \epsilon$$

$$P(\mu \in I_\epsilon) = 1 - \epsilon$$

$$P(\mu \in I_\epsilon^n) = (1-\epsilon)^n \xrightarrow{n \rightarrow \infty} 0$$



$$I_\epsilon \times I_\epsilon$$

$$\mu = (\mu_1, \mu_2)$$

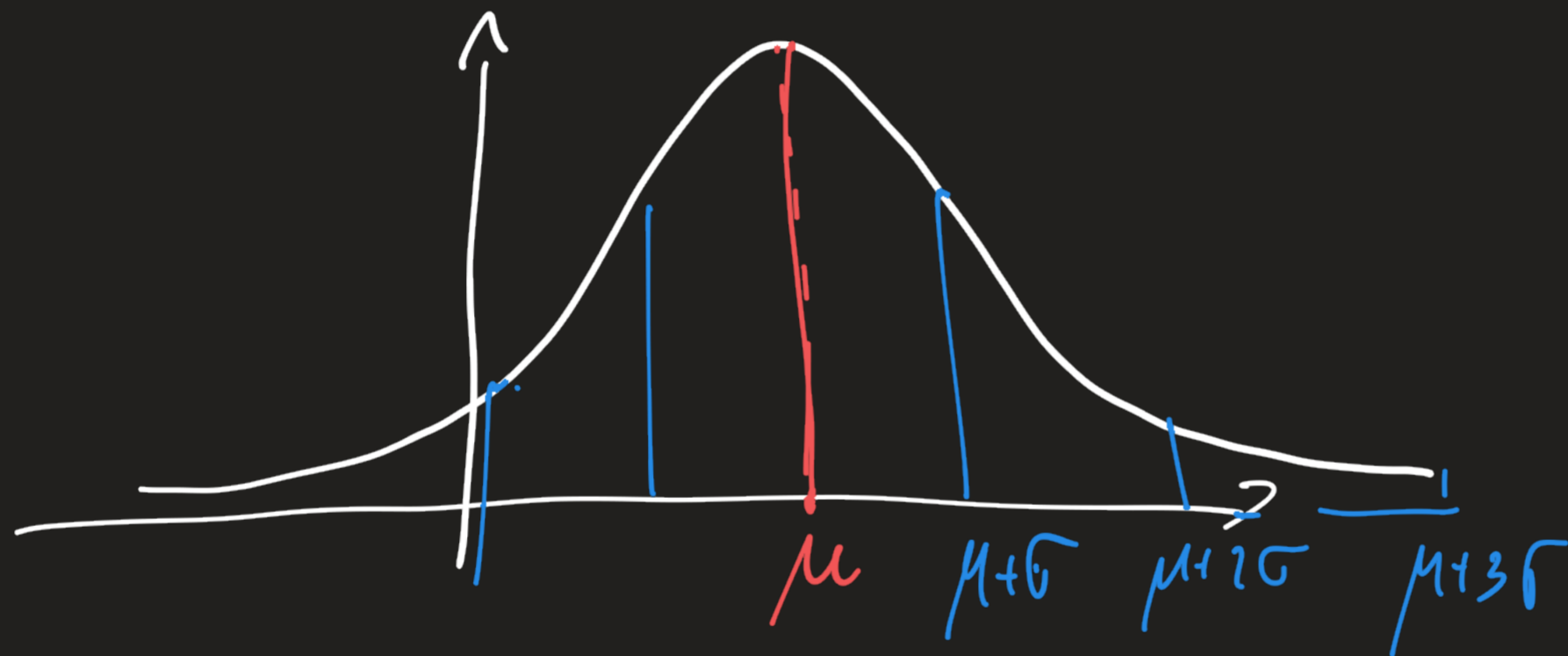
$$P(\mu \in I_\epsilon^2) = (1-\epsilon)^2$$

$\mu_j \sim U(0,1)$ ,  $\mu_j$  and  $\mu_i$  are indep.  $j \neq i$

# GAUSSIAN (NORMAL) DISTRIBUTION

$$P(x | \mu, \sigma^2) = \mathcal{N}(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

↑ mean      ↑ variance



$$P(\underline{x} | \underline{\mu}, \underline{\Sigma}) = \mathcal{N}(\underline{x} | \underline{\mu}, \underline{\Sigma}) = \frac{1}{\sqrt{(2\pi)^n \det(\underline{\Sigma})}} \exp\left(-\frac{1}{2} (\underline{x}-\underline{\mu})^T \underline{\Sigma}^{-1} (\underline{x}-\underline{\mu})\right)$$

↑ n-dim vector      ↑ n-dim vector      ↑

$\forall \underline{x}, \underline{x}^T \underline{\Sigma} \underline{x} \geq 0$

$\underline{\Sigma}_{ij} = \text{cov}(X_i, X_j)$   
 $i, j = 1, \dots, n$

# MIXTURES

$$p_1(x), \dots, p_k(x) \Rightarrow \text{prob. densities}$$

$x_1 \qquad \qquad \qquad x_k$

$$\pi_1, \dots, \pi_k, \sum \pi_i = 1, \pi_i \in [0, 1]$$

$$p(x) = \sum_{i=1}^k \pi_i p_i(x) \qquad p_i(x) = \mathcal{N}(x | \underline{\mu}_i, \underline{\Sigma}_i)$$

$\uparrow$

$$p(x, z) \neq p(x, z=i) = \underbrace{p(z=i)}_{\pi_i} \underbrace{p(x | z=i)}_{p_i(x)}$$

$z \in \{1, \dots, k\}$

$z \Rightarrow$  LATENT VARIABLE

$$p(x) = \sum_{i=1}^k p(x, z=i) = \sum_{i=1}^k p(z=i) p_i(x) = \sum_{i=1}^k \pi_i p_i(x)$$

$$= \sum_{z \in \mathcal{Z}} p(x, z)$$

$\nearrow$

## CONTINUOUS MIXTURES

$$P(x, z) = P(z) P(x|z)$$

$$P(x) = \int P(z) P(x|z) dz$$

$$P(z) = \text{GAMMA}(a, b)$$

$$P(x|z) = \mathcal{N}(x | \mu, z^{-1})$$

$$P(x) = t\text{-student}$$

with mean  $\mu$

$\nu = 2a$  degrees of freedom

$\lambda = \frac{a}{b}$  scale

# PROBABILISTIC INFERENCE

$H \in \{0, 1\}$  -  $H=1$  iff a person likes to work 14 hours per day

$F \in \{0, 1\}$  -  $F=1$  iff a person lives FRICO.

$$H=1 \underbrace{ \Rightarrow }_{\text{most likely}} F=1 \quad P(F=1 | H=1) = 0.8$$

$$H=1 \text{ is "rare"} \quad P(H=1) = 10^{-4}$$

• If  $F=1$  is common, how likely a  $F=1$  is also  $H=1$ .

$P(F=1) = 0.4$

$$P(H=1 | F=1) = \frac{P(F=1 | H=1) P(H=1)}{P(F=1)} = \frac{0.8 \cdot 10^{-4}}{0.4} = 2 \cdot 10^{-4}$$



$$P(F=1) = 2 \cdot 10^{-4}$$

$$P(H=1 | F=1) = \frac{P(F=1 | H=1) \cdot P(H=1)}{P(F=1)} = \frac{0.8 \cdot 10^{-4}}{2 \cdot 10^{-4}} = 0.4$$

$$P(F=1) = P(F=1 | H=1)P(H=1) + P(F=1 | H=0)P(H=0)$$

$$P(x|y) = \frac{P(y|x) \cdot P(x)}{P(y)}$$
$$\int P(y|x)P(x)dx$$

# MAXIMUM LIKELIHOOD

$X$ :  $x_1, \dots, x_N$  observation/data i.i.d

$P(X|\theta)$ , there is a  $\theta \in \Theta$  that generates the data \*

$$P(\underline{X}|\theta) = \prod_{i=1}^N P(x_i|\theta) \quad \text{AS A FUNCTION OF } \theta.$$

M.L. is FND  $\theta_{ML} = \arg \max_{\theta} P(\underline{X}|\theta)$

$$L(\theta) = \log P(\underline{X}|\theta) = \sum_{i=1}^N \log P(x_i|\theta)$$

$$\theta_{ML} = \arg \max_{\theta} L(\theta)$$

$$x_1, \dots, x_N \in \{0, 1\}$$

$$P(x|\theta) = \text{Bernoulli}(x|\theta) \quad , \theta \in [0, 1]$$
$$= \theta^x (1-\theta)^{1-x}$$

$$K = \sum_{i=1}^N x_i$$

$$L(\theta) = \sum_{i=1}^N \log P(x_i|\theta) = \sum_{i=1}^N x_i \cdot \log(\theta) + (1-x_i) \log(1-\theta)$$
$$= K \cdot \log \theta + (N-K) \log(1-\theta)$$

$$0: \frac{\partial L(\theta)}{\partial \theta} = \frac{K}{\theta} - \frac{(N-K)}{1-\theta} = 0 \Rightarrow \theta = \frac{K}{N}$$

# MAXIMUM A-POSTERIORI

$\underline{x} = x_1 \dots x_n$  iid

$p(x|\theta)$  - model

$p(\theta)$  - prior

$$\underbrace{p(\theta|\underline{x})}_{\text{POSTERIOR ON } \theta \text{ GIVEN DATA}} = \frac{p(\underline{x}|\theta) \cdot p(\theta)}{\underbrace{p(\underline{x})}_{\text{MARGINAL LIKELIHOOD}}}$$

$$\theta_{\text{MAP}} = \underset{\theta}{\text{argmax}} p(\theta|\underline{x})$$

$$= \underset{\theta}{\text{argmax}} \log p(\underline{x}|\theta) + \log p(\theta)$$

# BAYESIAN ESTIMATION

$\underline{x}: x_1 \rightarrow x_N$  iid

$p(x|\theta)$  LIKELIHOOD

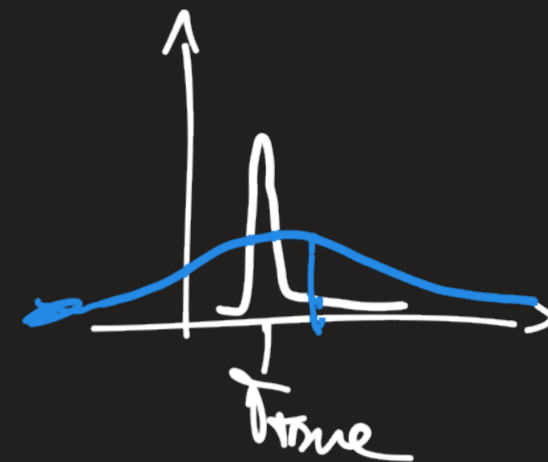
$p(\theta)$  PRIOR

$$p(\theta|\underline{x}) = \frac{p(\underline{x}|\theta) p(\theta)}{p(\underline{x})}$$

✓  $p(\underline{x})$

$$p(\underline{x}|\underline{x}) = \int p(x|\theta) p(\theta|\underline{x}) d\theta$$

$$\int p(x, \theta|\underline{x}) d\theta = \int p(x|\theta) p(\theta|\underline{x}) d\theta$$



POSTERIOR

$$p(\underline{x}) = \int p(\underline{x}|\theta) p(\theta) d\theta$$

PREDICTIVE DISTRIBUTION

$$x_1 \dots x_N \in \{0,1\} \quad p(x|\theta) = \text{Bernoulli}(\theta) \quad \theta \in [0,1]$$

$$p(\theta) = \text{Beta}(\theta|\alpha, \beta) = \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

$$\mathbb{E}_{\text{Beta}(\alpha, \beta)}[\theta] = \frac{\alpha}{\alpha + \beta} \quad k = \sum x_i$$

$$\begin{aligned} \log p(\theta|x) &= L(\theta) + \log B(\theta|\alpha, \beta) - \log p(x) \\ &= C + k \log \theta + (N-k) \log(1-\theta) + (\alpha-1) \log \theta + (\beta-1) \log(1-\theta) \\ &= C + \underbrace{(k+\alpha-1) \log \theta + (N-k+\beta-1) \log(1-\theta)} \\ &= \log \text{Beta}(\theta|k+\alpha, N-k+\beta) \quad p(x=1|x) = \end{aligned}$$

Beta is Bernoulli  
CONJUGATE PRIOR

$$= \frac{k+\alpha}{N+\alpha+\beta}$$

# INFORMATION THEORY PILLS - ENTROPY

$p(x)$

$-\log p(x)$  - SELF INFORMATION

$$H[p] = \mathbb{E}_p[-\log p(x)] = - \int p(x) \log p(x) dx$$

max is  
Normal R.V.

$$H[p] = - \sum_i p(x_i) \log p(x_i) \leftarrow \begin{array}{l} \text{max for unif. distr.} \\ \log K. \end{array}$$

# KULLBACK LEIBLER DIVERGENCE

RELATIVE ENTROPY  $P, Q$

$$KL[Q \parallel P] = \int Q(x) \log \frac{Q(x)}{P(x)} dx = -H[Q] - \underbrace{E_Q[\log P]}$$

$$KL[Q \parallel P] = 0 \iff Q = P$$

KL is a convex functional and  $KL \geq 0$

$\rightarrow$   $P$  is fixed but unknown  $Q = Q_\theta$  can vary  
What is the best  $Q_\theta$  that approx  $P$ .

$\theta^* = \underset{\theta}{\operatorname{argmin}} KL[Q_\theta \parallel P]$  } AN VARIATIONAL INFERENCE

$I[X, Y] = KL[P(X, Y) \parallel P(X)P(Y)]$  MUTUAL INFORMATION



$$\underline{x} = x_1, \dots, x_N$$

EMPIRICAL DISTRIBUTION  $P_{emp}(x) = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(x=x_i)$

$q$

$$K_L [P_{emp} \parallel q] = \underbrace{E_{P_{emp}}[-\log q(x)]}_{\substack{\uparrow \\ -\frac{1}{N} \sum \log q(x_i)}}} - H[P_{emp}]$$

$-\frac{1}{N} \sum \log q(x_i)$  CROSS ENTROPY

$\forall q = q_\theta$

$$-\frac{1}{N} L''(\theta)$$



MAXIMIZING  $L(\theta)$



MINIMIZING K.L. divergence between  $P_{emp}$  and  $q_\theta$ .