

Stats/ Probability

"E" a random event.

S: set of all possible outcomes associated with E.  
Sample space

E.g. E = coin flip -> S = {H, T}

Random Variable: X (a R.V.) is assigned

a number according to outcome of a random event.

R.V.'s: Discrete or Continuous  
e.g. Yes/No, coin flip      e.g. height, weight, time, etc.

Probability Distributions  
called a PMF (point-mass) for Discrete R.V.'s  
called a density for Contin. R.V.'s.

(1)  $0 \leq P(X=i) \leq 1$  for all  $i \in S$

(2)  $\sum_{i \in S} P(X=i) = 1$        $\left( \int_{-\infty}^{\infty} P(x) dx = 1 \right)$



(Note:  $P(x_i) = P(X = i)$ , etc.)

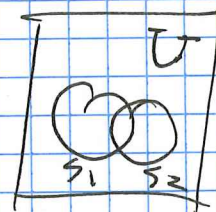
11

Events in  $S$  are disjoint if:  $S_1 \cap S_2 = \emptyset$   
( $S_1, S_2 \in S$ )

If  $S_1, S_2$  disjoint:

Then:  $P(S_1 \text{ OR } S_2) = P(S_1) + P(S_2)$

More Generally, Additive Rule of Prob:

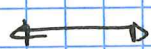


$$P(S_1 \text{ OR } S_2) = P(S_1) + P(S_2) - P(S_1 \text{ AND } S_2)$$

Conditional Probability

$$P(A|B) \stackrel{\text{def.}}{=} \frac{P(A \& B)}{P(B)}$$

Prob. of A, given B



$$P(A \& B) = P(A|B) \cdot P(B)$$

Multiplication Rule

$$P(A \& B) = P(A|B) \cdot P(B) = P(B|A) \cdot P(A)$$



# Independence

We say events  $A$  &  $B$  are independent if outcome of  $A$  has no bearing on  $B$  & vice versa.

IF  $A$  &  $B$  are independent more formally:

$$P(A \& B) = P(A)P(B) \quad (\text{ie. Re joint probability factors})$$

Also, equivalently:

IF  $A, B$  independent then:

$$P(A|B) = P(A), \quad P(B|A) = P(B)$$

Thus, if  $A, B$  independent:

$$P(A \& B) = \underbrace{P(A|B)}_{B, \text{ Mult. Rule}} \cdot P(B) = \underbrace{P(A)}_{B, \text{ independence}} \cdot P(B)$$

equivalence



② Major Theorems in Elementary STATS:

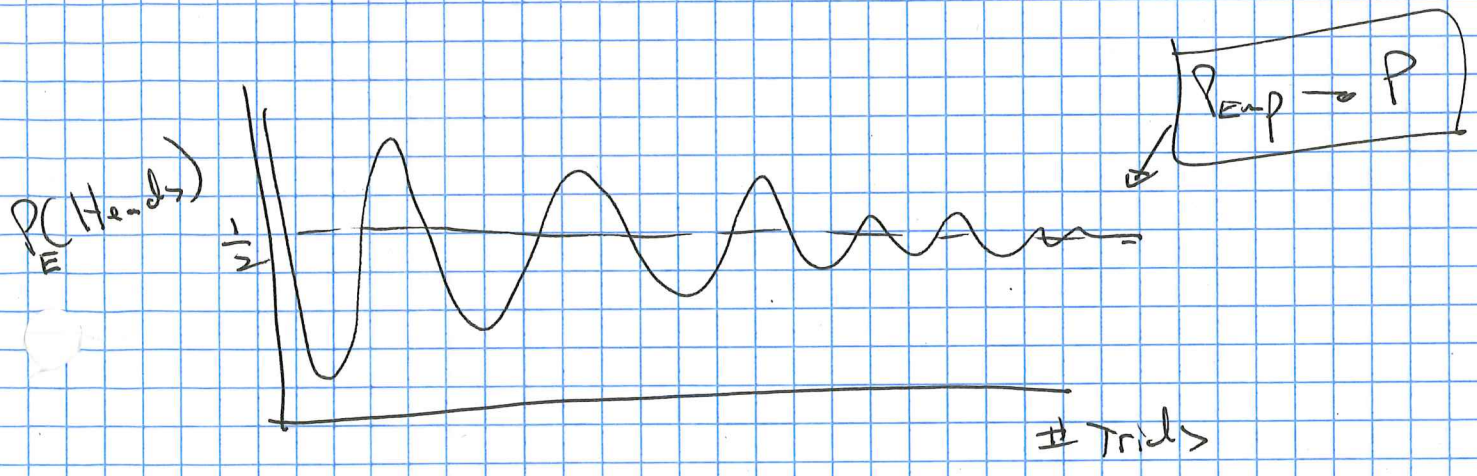
- ① Law of Large Numbers (LLN)
- ② Central Limit Theorem (CLT)

**LLN**: (Paraphrasing) Experimental (i.e. empirical probabilities) <sup>Probs</sup> converge to their associated Theoretical probability density as the number of Trials tends to infinity.

$$\lim_{n \rightarrow \infty} P_{\text{Exp}}(x) = P(x)$$

e.g. Consider a single flip of a fair coin, i.e.

$$P(\text{Heads}) = \frac{1}{2} \text{ ("True" or Theoretical probability)}$$





**CDF** (Cumulative Density Function)

$$F_X(x) = P(X \leq x) = \int_{-\infty}^x P(u) du$$

$F_X(x) = .25$   
( $x = Q_1$  Quantile)

$F_X(x) = .5$   
 $x = Q_2$ : Median

$F_X(x) = .75$   
 $x = Q_3$

Note:  $\frac{d}{dx} F(x) = P(x)$  why?  
cdf                      pdf

Some Essential Distributions

Normal / Gaussian  
(Continuous)

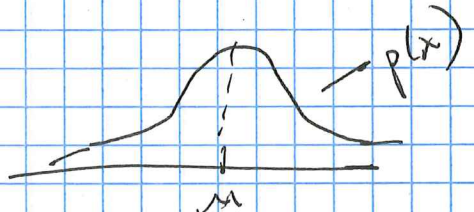
(1-D)

we write:  
 $X \sim N(\mu, \sigma)$

$$P(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$$

Normalized constant

Standard Normal:  $N(0, 1)$



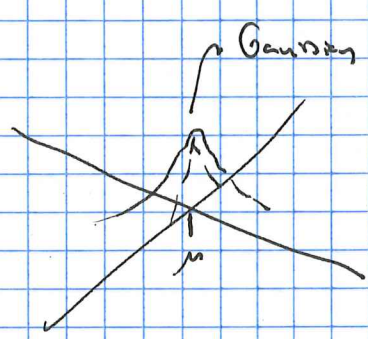
Empirical Rule

- $x \sim N(\mu, \sigma)$
- $x \pm \mu$ : 68%
- $x \pm 2\mu$ : 95%
- $x \pm 3\mu$ : 99.7%



MVN (Multi-variate Normal)

P(x) = 1 / (sqrt(2pi)^d |Sigma|) exp[-1/2 (x - mu)^T Sigma^-1 (x - mu)]



Note: Delta = Sigma^-1

referred to as the precision matrix

Phi(x) = integral from -infinity to x of P(z) dz = 1/2 (1 + erf((x - mu) / (sigma \* sqrt(2))))

where: erf(.) is the error function, has no closed-form expression.

Bernoulli (Discrete)

P(X=1) = theta
P(X=0) = 1 - theta

0 <= theta <= 1

e.g. flip a coin w/ Prob. Heads = theta
Tails = 1 - theta



Binomial Distribution

(Discrete)

Binary outcome

$$S = \{0, 1\}$$

(.) Case of n Bernoulli trials

let:  $p(X=1) = \theta$  "success"  
 $p(X=0) = 1-\theta$  "failure"

$$P(X=k) = \binom{n}{k} \theta^k (1-\theta)^{n-k}$$

prob. "k successes in n trials"

ex. flip biased coin 10 times  $p(H) = .6$   
 $p(T) = .4$

$$P(\text{exactly 7 H in 10 flips}) = \binom{10}{7} (.6)^7 (.4)^3$$

Poisson Distribution

(Discrete)

$\lambda$ : Mean # successes in given time interval

$$P(X=k) = \frac{e^{-\lambda} \lambda^k}{k!}$$

$$k = 0, 1, 2, \dots$$

Q: If, on average, 4 people visit a given webpage per minute, what is prob. of two or fewer visitors?



$P(\text{Two or fewer visitors in 4 min}) =$

$$P(X=0) + P(X=1) + P(X=2)$$

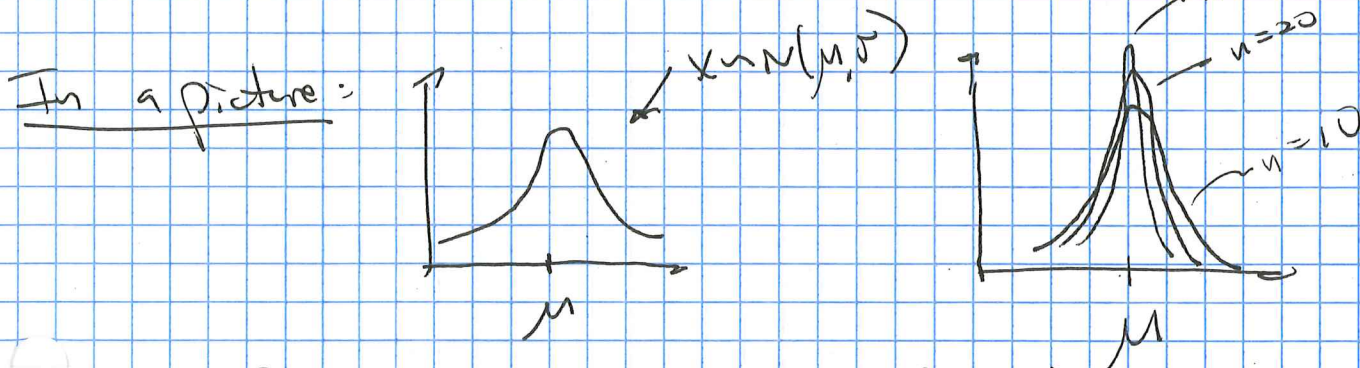
$$= \frac{e^{-4} \cdot 4^0}{0!} + \frac{e^{-4} \cdot 4^1}{1!} + \frac{e^{-4} \cdot 4^2}{2!} \approx \boxed{.238}$$

**CLT** (Central Limit Theorem)  
(Classical, non-technical version)

Given:  $x_1, x_2, \dots, x_n$  **IID**  $\rightarrow$  Independent, Identically Distributed  
Random sample

where:  $x_i \sim N(\mu, \sigma^2)$

Then:  $\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$



Q: Given IQ scores have dist:  $N(100, 15)$

what is prob:  $P(85 \leq X \leq 115) = 68\%$

For 10 individual, what is:  $P(85 \leq \bar{X} \leq 115) = 99\%$



# Expectation (of a RV)

$$E[X] = \sum_i x_i P(X=x_i)$$

(Discrete case)

$$E[X] = \int_{-\infty}^{\infty} x P(x) dx$$

(Continuous)

Q: How many Heads are expected in 10 flips of a fair coin?

X	P(X)
0 ( $\bar{H}$ )	$\frac{1}{2}$
1 (H)	$\frac{1}{2}$

→ PMF (Bernoulli Trials)  
(n=10)

$$\begin{aligned}
 E[X] &= \sum_{k=0}^{10} \binom{10}{k} \left(\frac{1}{2}\right)^k \left(\frac{1}{2}\right)^{10-k} \cdot k \\
 &= \sum_{k=0}^{10} \binom{10}{k} \left(\frac{1}{2}\right)^{10} \cdot k = \left(\frac{1}{2}\right)^{10} \sum_{k=0}^{10} \binom{10}{k} \cdot k \\
 &= \boxed{5}
 \end{aligned}$$



# Variance & Standard Deviation (for RVs)

19

$$\text{Var}[X] = E[(X - \mu)^2] = \sum_i (x_i - E[X])^2 (P(X=x_i))$$

$\mu = E[X]$

$$\text{SD}[X] = \sqrt{\text{Var}[X]}$$

Corollary:  $\text{Var}[X] = E[X^2] - \mu^2$

Proof:

$$\text{Var}[X] = E[(X - \mu)^2] = E[X^2 - 2\mu X + \mu^2]$$

$$= E[X^2] - E[2\mu X] + E[\mu^2]$$

By "linearity" of Expectation

$\mu^2$ : a constant

$$= E[X^2] - 2\mu \cdot E[X] + \mu^2$$

by linearity

$E[\text{constant}] = \text{constant}$

$$= E[X^2] - 2\mu \cdot \mu + \mu^2 = E[X^2] - \mu^2 \quad \square$$



# Covariance

$X, Y$  R.V.'s:

$$\text{Cov}(X, Y) = E[(X - \mu_x)(Y - \mu_y)]$$

Lemma: If  $X, Y$  are independent, then  $\text{Cov}(X, Y) = 0$ .

Pf:  $E[(X - \mu_x)(Y - \mu_y)] =$

$$E[X Y - X \mu_y - Y \mu_x + \mu_x \mu_y] \stackrel{?}{=}$$

$$\stackrel{?}{=} E[X Y] - \mu_y E[X] - \mu_x E[Y] + E[\mu_x \mu_y]$$

By  
Linearity of  
E[.]

$$= E[X Y] - \mu_y \mu_x - \mu_x \mu_y + \mu_x \mu_y$$

$$\left( E[X Y] = \sum_k \sum_i p(x, y) \cdot x y = \sum_k \sum_i p(x) \cdot x \cdot p(y) \cdot y = E[X] \cdot E[Y] \right)$$

$$= E[X] \cdot E[Y] - \mu_x \mu_y$$

$$= \mu_x \cdot \mu_y - \mu_x \cdot \mu_y = 0. \quad \square$$



## Covariance Matrix

(2)

let  $\vec{X} = (x_1, \dots, x_n)$  (a vector of RV's)

$$\Sigma_{ij} = \text{Cov}(X_i, X_j) = E[(X_i - \mu_i)(X_j - \mu_j)]$$

Matrix  
of  
Covariances

Note:  $\Sigma_{ij} =$

$$\begin{bmatrix} \text{Var}[X_1] & \text{Cov}(X_1, X_2) & \dots & \text{Cov}(X_1, X_n) \\ | & | & & | \\ | & | & & | \\ | & | & & | \\ \text{Cov}(X_n, X_1) & \dots & & \text{Var}[X_n] \end{bmatrix}$$

$\Sigma$  is symmetric, positive, semi-definite.

## Bayes' Theorem

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

↳ Re "Bayes' Theorem" of AI/ML.

Pf:

$$P(A|B) = \frac{P(A \& B)}{P(B)} = \frac{P(B|A) \cdot P(A)}{P(B)} \quad \square$$



## Medical Diagnosis Example

22

Q: Physician knows  $P(\text{Disease}) = 1\%$   
 $P(\text{No Disease}) = 99\%$

$$P(\oplus \text{ Test} | \text{D}) = 0.792 \quad (\text{True Positive})$$

$$P(\oplus \text{ Test} | \text{No D}) = 0.096 \quad (\text{False Positive})$$

Given That Patient ~~X~~ receives  $\oplus$  Test result,

what is probability they have disease?

Want:  $P(\text{D} | \oplus) = \frac{P(\oplus | \text{D}) P(\text{D})}{P(\oplus)}$

By Bayes'

$$= \frac{P(\oplus | \text{D}) P(\text{D})}{P(\oplus | \text{D}) P(\text{D}) + P(\oplus | \text{No D}) P(\text{No D})} = \frac{(0.792)(0.01)}{(0.792)(0.01) + (0.096)(0.99)}$$

$$\approx \boxed{0.103}$$



# Frequentist vs. Bayesian Statistics

Frequentists = Model  $\theta$  parameters are fixed (i.e. Platonic);  
Data are drawn from "God's distribution", defined by  $\theta$ .

Bayesian: Data are fixed (the observed data);  
data are observed from realized sample; we encode prior beliefs; parameters are described probabilistically.

Frequentists: Use MLE (Maximum likelihood estimate) for point estimate of  $\theta$ :  
("likelihood" of data)

$$\hat{\theta}_{MLE} = \arg \max_{\theta} L(D|\theta)$$

Bayesian: Compute MAP (Maximum a Posterior) estimate:  
("prior" on  $\theta$ )

$$\hat{\theta}_{MAP} = \arg \max_{\theta} P(D|\theta)P(\theta)$$

where posterior =  $p(\theta|D) \propto p(D|\theta)p(\theta)$



Various "Pathologies" of Frequentist Prob. Exist:

The Problem of Induction (Hume), Black Swan Paradox, limited exact solutions, reliance on "long-term" frequencies.

Bayesian Statistics have been called "Statistics of the 21st cent." (Efron)

Q: Do I need any Calculus for AI/ML?

A: Honestly, not that much!

At minimum, know 2 Basic Things:

1) Calculus gives us a framework to coherently/mathematically describe the flow/dynamics of the real-world. (engineers, physicists, etc. like this)

2) Principles in Calculus are helpful in optimization!

If  $f(x)$  is convex / concave

it has a unique, global minimum/maximum.

We can use differential calculus to solve convex/concave problems (or even approximate sol. for non-convex)



# Hill-Climbing

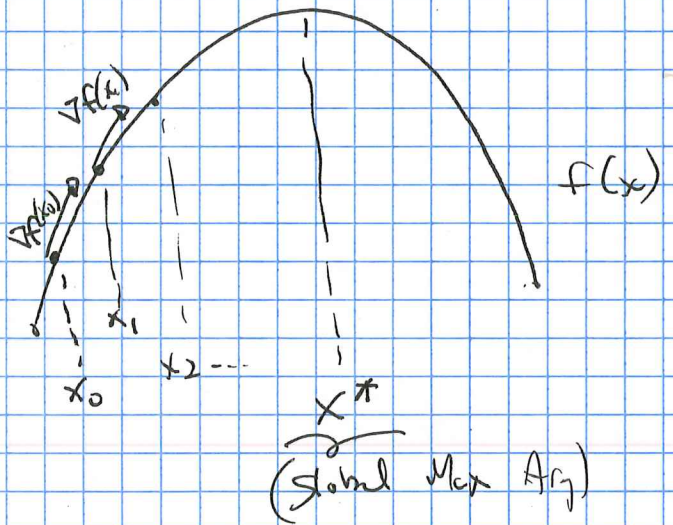
Suppose  $f(x)$  is concave - we devise an iterative,  
 "hill-climbing" algorithm to approximate:  $x^* = \arg \max_x f(x)$

(1) Initial guess:  $x_0$

(2) Compute gradient @  $x_0$ :

$$\nabla f(x_0)$$

Recall:  $\nabla f = \langle f_{x_1}, f_{x_2}, \dots, f_{x_n} \rangle$   
 vector of partial derivatives



(3) Set  $x_1 = x_0 + \delta \nabla f(x_0)$

$\delta$ : learning parameter  
 (controls convergence speed)

loop: (2)-(3) for  $x_2, x_3, \dots$

until stopping condition.



# (Very Brief) Information Theory

## Entropy of a RV

$$H(X) = - \sum_i p(x=i) \log_2 p(x=i)$$

(Define:  $0 \log 0 \equiv 0$ )

quantifies disorder/uncertainty

eg. By some accounts, entropy for English language (per letter)

is  $\approx 2.62$  bits. (using N-gram model), i.e.

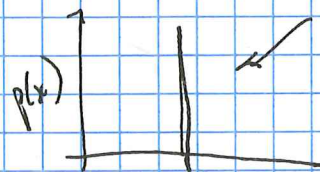
on average we need  $\approx 2.62$  binary questions to guess

$n$ -th letter of a string. (Also a measure of redundancy)

(\*) The distribution of Max entropy is the uniform distribution.  $\rightarrow$  

(\*) The distribution of Minimum entropy is the

(Pinc) delta function.  $\rightarrow$



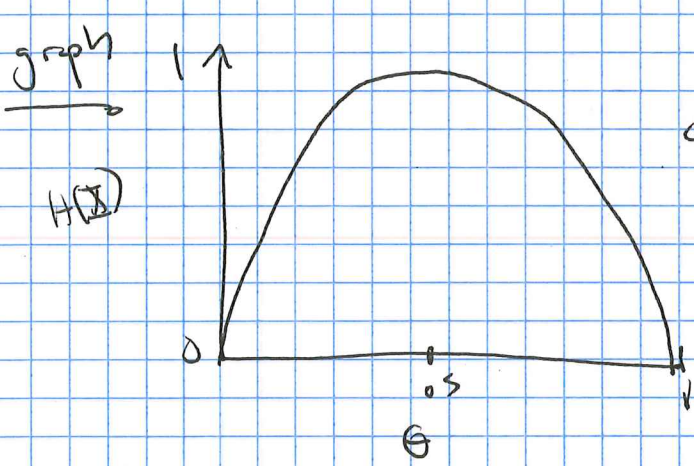
deterministic,  
i.e. zero  
uncertainty



Ex. Consider Bernoulli RV:  $X$

$$P(X=1) = \theta, \quad P(X=0) = 1-\theta$$

$$\begin{aligned}
 H(X) &= - \sum_x p(x) \log_2 p(x) \\
 &= - \left[ p(X=1) \log_2 p(X=1) + p(X=0) \log_2 p(X=0) \right] \\
 &= - \left[ \theta \log_2 \theta + (1-\theta) \log_2 (1-\theta) \right]
 \end{aligned}$$



Max occurs when  $\theta = 0.5$  (i.e. uniform);

Min occurs when  $\theta = 0$  or  $1$  (i.e. deterministic)

**KL Divergence** (Kullback-Leibler)

( $p, q$  are prob. distributions)

$$KL(p \parallel q) = \sum_i p_i \log \frac{p_i}{q_i}$$

Measure of dissimilarity b/w distributions

$$\begin{aligned}
 &= \sum_i p_i \log p_i - \sum_i p_i \log q_i \\
 &= -H(p) + H(p, q)
 \end{aligned}$$

entropy      "cross entropy"



Note:  $KL(p||q) \geq 0 \neq KL(p||q) = 0 \text{ iff } p=q.$

"Information Inequality"

Recall: Covariance (& correlation) measure the linear dependence b/w RVs.

Using KL-Divergence we can develop a more general

notion of dependence: Mutual Information (MI)

$$I(X; Y) = KL(p(x, Y) || p(x)p(Y))$$

$$= \sum_x \sum_y p(x, y) \log(p(x)p(y))$$

From above:  $I(X; Y) \geq 0 \neq I(X; Y) = 0 \text{ iff } p(x, y) = p(x)p(y)$

Thus:  $I(X; Y)$  measures "similarity" between  $p(x, y)$  &  $p(x)p(y)$

$p(x)p(y)$   
 $\xrightarrow{\text{factored}}$   
 joint

Shannon Source Coding Theorem:

$$H(X) \leq E[L] \leq H(X) + 1$$

(for binary alphabet)

$E[L]$ : expected word length