

\rightarrow Prob. b. lity

" E " a random event.

S : set of all possible outcomes associated w.r.t E .
 Sample space

e.g. $E = \text{Coin flip} \rightarrow S = \{H, T\}$.

Random Variable : X (\sim R.V.) \mapsto assigned

a number according to outcome of a random event.

R.V's: Discrete

e.g. Yes/No,
coinflip

or Continuous

e.g. height, weight, Time, etc.

Probability Distributions

→ called \rightarrow (PMF) (point-mass)
for Discrete R.V's

→ called \rightarrow (Density) for Continuous R.V's.

$$(1) 0 \leq P(X=i) \leq 1 \quad \text{for all } i \in S$$

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

(Note: $P(x_i) = P(X = x_i, \omega_{\text{rc}})$)

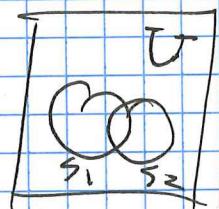
11

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

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

Prob. of A, given B

$$\leftrightarrow (P(A \text{ and } B) = P(A|B) \cdot P(B))$$

Multiplication Rule

:

$$P(A \text{ and } B) = P(A|B) \cdot P(B) = P(B|A) \cdot P(A)$$

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:

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

(i.e. 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:

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

\curvearrowleft Mult Rule

\curvearrowright A, B independent

equivalence

(2) Major Theorems in Elementary Stats:

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

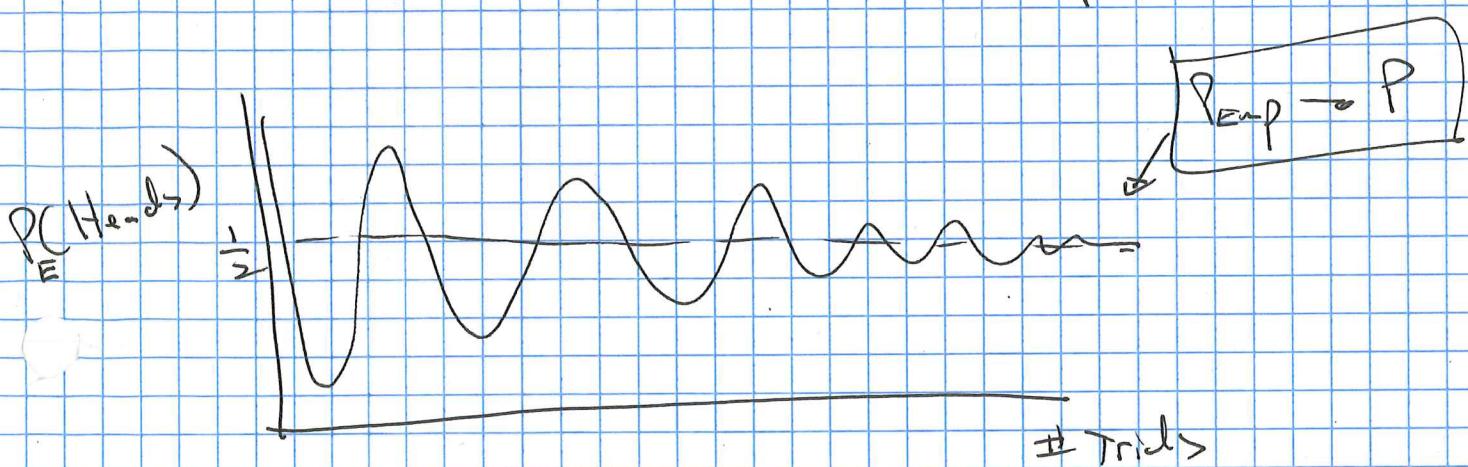
LLN: (Paraphrasing) Experimental (i.e. empirical

probabilities) ^{probs} converge To Their associated Theoretical
probability density as The number of Trials Tends
To infinity.

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

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

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



CDF (Cumulative Density Function)

(14)

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

$$F_X(x) = .25$$

$x = Q_1$ Quantile

$$F_X(x) = .5$$

$x = Q_2$: Median

$$F_X(x) = .75$$

$x = Q_3$

Note: $\frac{d}{dx} F(x) = P(x)$

$\frac{\partial}{\partial x}$ cdf $\frac{\partial}{\partial x}$ pdf

why?

Some Essential Distributions

Normal / Gaussian
(continuous)

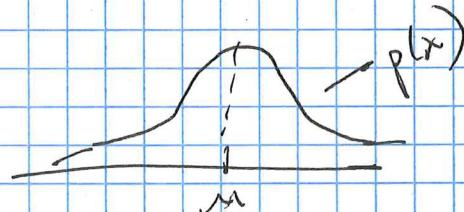
(1-D)

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

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

$\frac{1}{\sigma \sqrt{2\pi}}$ Normalization constant

Standard Normal: $N(0, 1)$



Empirical Rule

$x \sim N(\mu, \sigma)$

$x \pm \mu : 68\%$

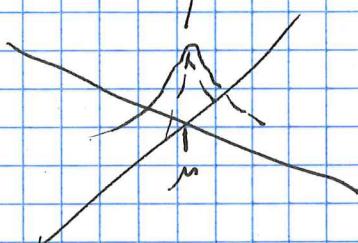
$x \pm 2\mu : 95\%$

$x \pm 3\mu : 99.7\%$

MVN (Multi-variate Normal)

$$P(\vec{x}) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp \left[-\frac{1}{2} (\vec{x} - \vec{\mu})^\top \Sigma^{-1} (\vec{x} - \vec{\mu}) \right] \quad (\vec{x} \in \mathbb{R}^d)$$

Gaussian "Mound" in \mathbb{R}^2



Note: $\Delta = \Sigma^{-1}$

referred to as Precision matrix

$$\Phi(x) = \int_{-\infty}^x P(z) dz = \frac{1}{2} \left(1 + \frac{\text{erf}(z - \mu)}{\sigma \sqrt{2}} \right)$$

(Cumulative
Normal)

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

Bernoulli:
(Discrete)

$$P(X=1) = \theta$$

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

$$0 \leq \theta \leq 1$$

E.g. flip a coin w/ Prob. Heads = θ
Tails = $1 - \theta$

16

Binomial Distribution
(Discrete)

Binary Outcome

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

(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 $P(H) = .6$
 $P(T) = .4$

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

Poisson Distribution
(Discrete)

λ : 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?

17

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

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

$$= \frac{e^{-4} \cdot 4^0}{0!} + \frac{e^{-4} \cdot 4^1}{1!} + \frac{e^{-4} \cdot 4^2}{2!} \approx 0.238$$

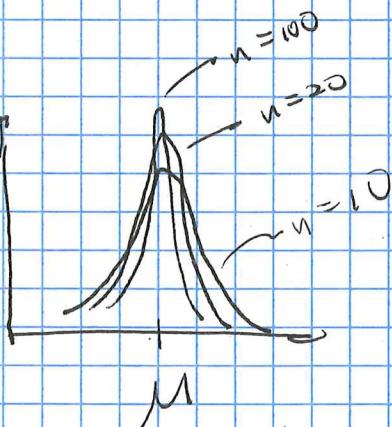
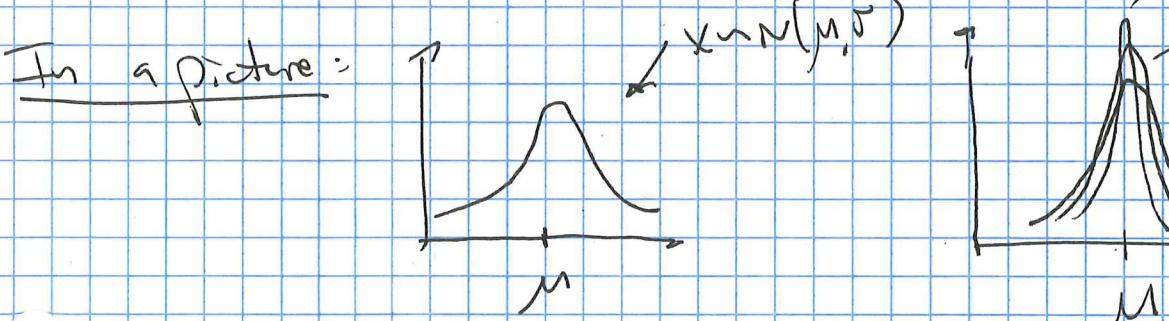
CLT Central Limit Theorem
(Classical, non-Technical version)

Given: x_1, x_2, \dots, x_n
 Random sample

IID → Independent, Identically
Distributed

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

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



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

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

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

18

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 flip of a fair coin?

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

RMF (Bernoulli Trials)
(n=10)

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

$$= \sum_{k=0}^{10} \binom{10}{k} \left(\frac{1}{2}\right)^{10} \cdot k = \left(\frac{1}{2}\right)^{10} \sum_{k=0}^{10} \binom{10}{k} \cdot k$$

$$= [5]$$

19

Variance & Stand Deviation (for RVs)

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

$\mu = E[X]$

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

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

Proof:

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

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

↓ μ^2 : a constant

By "linearity" of Expectation

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

↑ linearity

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

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

Covariance

X, Y Rv's:

$$\boxed{\text{Cov}(X, Y) = E[(X - \mu_X)(Y - \mu_Y)]}$$

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

Pf: $E[(X - \mu_X)(Y - \mu_Y)] =$

$$E[XY - X\mu_Y - Y\mu_X + \mu_X\mu_Y] \ni$$

$$\ni E[XY] - \mu_Y E[X] - \mu_X E[Y] + E[\mu_X\mu_Y]$$

By
Linearity of

$$E[\cdot] = \underbrace{E[XY]}_{=} - \mu_Y \mu_X - \cancel{\mu_X \mu_Y} + \cancel{\mu_X \mu_Y}$$

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

since X, Y independent

$$= E[X] \cdot E[Y] - \mu_X \mu_Y$$

$$= \mu_X \cdot \mu_Y - \mu_X \cdot \mu_Y = 0.$$

□

21

Covariance Matrix

Let $\bar{X} = \langle X_1, \dots, X_n \rangle$ (a vector of RV's)

$$\sum_{ij} = \text{Cov}(X_i, X_j) = E[(\bar{X}_i - \mu_i)(\bar{X}_j - \mu_j)]$$

Matrix
of
(covars)

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

\sum is symmetric, positive, semi-definite.

Bayes' Theorem

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

L Re "Bayes Theorem" of AI/ML.

Pf:

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

(Medical Diagnosis Example)

(22)

Q: Physician knows $P(D_{\text{true}}) = 1\%$

$$P(\text{No } D_{\text{true}}) = 99\%$$

$$P(\text{+ Test} | D) = .792 \quad (\text{True Positive})$$

$$P(\text{+ Test} | \text{No } D) = .096 \quad (\text{False Positive})$$

Given That Patient \neq receives (+) Test result,

what is probability they have disease?

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

By Bayes'

$$= \frac{P(+|D)P(D)}{P(+|D)P(D) + P(+|\text{No } D)P(\text{No } D)}$$

$$= \frac{(.792)(.01)}{(.792)(.01) + (.096)(.99)}$$

$$\approx \boxed{.0103}$$

Frequentist vs. Bayesian Statistics

Frequentists \Rightarrow Model parameters are fixed (i.e. Picardie);

Data are drawn from "God's distribution", defined by θ .

Bayesians: 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 θ :

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

Bayesians: Compute MAP (Maximum Posterior) estimate:

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

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

Variety "Pathologies" of Frequentist Prob. Exist:

The Problem of Induction (Hume), Black Swan Paradox,
Limited exact solution, 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/dynamism 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 size for non-convex)

Hill Climbing

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

(1) Initial guess: x_0

(2) Compute gradient $\nabla f(x_0)$

$\boxed{\nabla f(x_0)}$

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

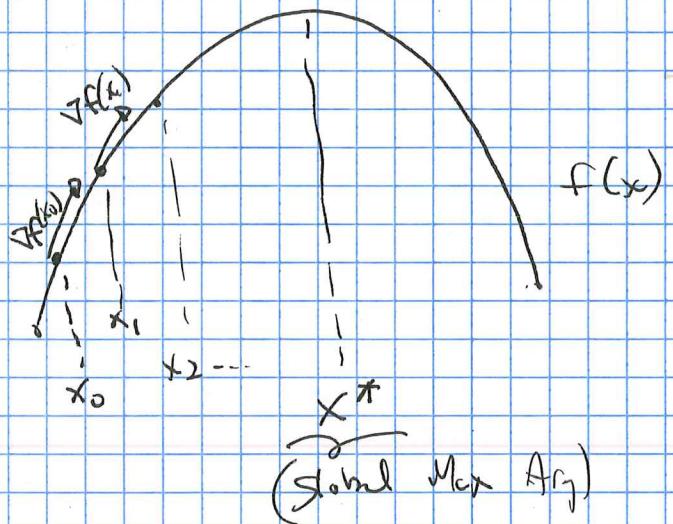
(3) Set $\boxed{x_1 = x_0 + \delta \nabla f(x_0)}$

δ : 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: $\log_0 \equiv 0$)

\downarrow
quantifies disorder/uncertainty

e.g. By some accounts, entropy for English language (per-letter)

is ≈ 2.62 bits, (using N-gram model), i.e.

on average we need ≈ 2.62 binary questions to guess

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

(*) The distribution of Max entropy \Rightarrow The uniform distribution. $\rightarrow p(x) \downarrow$

(*) The distribution of Minimum entropy is The

(Pmly) delta function. $\rightarrow p(x) \downarrow$ Deterministic, if zero uncertainty

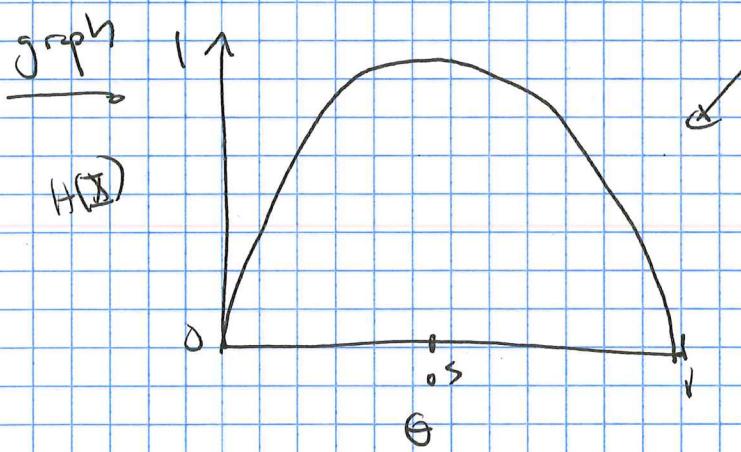
[Ex.] Consider Bernoulli RV: X

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

$$H(X) = - \sum p(x) \log_2 p(x)$$

$$= -[P(X=1) \log_2 P(X=1) + P(X=0) \log_2 P(X=0)]$$

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



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)

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

Measure of dissimilarity
b/w distributions.

"entropy"
"cross-entropy"

$$= \sum_i p_i \log p_i - \sum_i p_i \log q_i$$

$$= -H(p) + H(p, q)$$

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

"Information Inequality"

Recall: Covariance (& Correlation) measure the linear dependence b/w R.S.

Using KL-Divergence we can develop a more general notion of dependence: [Mutual Information (MI)]

$$I(X;Y) = \text{KL}(p(X,Y) || p(X)p(Y))$$

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

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

Thus: $I(X;Y)$ measures "similarity" between $p(x,y)$ & $\underbrace{p(x)p(y)}_{\text{joint}}$

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

(for binary alphabet)

$E[L]$: expected word length

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

Shannon Source Coding Theorem: