

Physics 509 Theory of Measurements Course Notes

Tobias Faehndrich

This document was last typeset on November 27, 2025

Introduction:

Notes written at UBC 2025W1 with Dr. Colin Gay.

Contents

1 Tuesday, September 9th 2025	6
1.1 Motivation: Stochastic Nature of Experimental Data	6
1.2 Probabilistic Interpretation of Experimental Results	6
1.3 Sample Spaces and Stochastic Variables	6
1.4 Events and Set Operations	7
1.5 Kolmogorov's Axioms of Probability	7
1.6 Consequences of the Probability Axioms	8
1.7 Uniform Probability on Finite Sample Spaces	8
1.8 Example: Probability of a Straight in Poker	9
1.9 Conditional Probability	9
2 Thursday, September 11th 2025	10
2.1 Bayes' Formula	10
2.2 Law of Total Probability	10
2.3 Independent Events	10
2.3.1 Example: Rolling Two Dice	11
2.4 Random Variables and Probability Distributions	11
2.5 Describing a Distribution	12
2.6 Cumulative Distribution Function (CDF)	12
2.7 Expectation Values	12
2.8 Characteristic Function	12
2.9 Central Moments	13
3 Tuesday, September 16th 2025	14
3.1 Bayes Theorem and Its Applications	14
3.2 The Monty Hall Problem: A Bayesian Analysis	14
3.3 Alternate Monty Hall Formulations	15
3.4 Monty Hall Generalized to n Doors	15
3.5 Continuous Probability Distributions and Moments	16
3.6 Discrete Probability Distributions	16
3.7 Cumulative Distribution Functions	17
3.8 Multivariate Distributions and Covariance	17
4 Thursday, September 18th 2025	18
4.1 Conditional Probability: A Simple Example	18
4.2 Distributions of Multiple Random Variables	18
4.3 Covariance Matrix and Correlation Coefficient	18
4.4 Independence and Uncorrelated Variables	19
4.5 Examples of Correlated and Uncorrelated Variables	19

4.6	Marginal Distributions	19
4.7	Change of Variables in Probability Densities	20
4.8	Change of Variables: Non One-to-One Case	20
4.9	Multivariate Transformations and the Jacobian	21
4.10	Example: Cartesian to Polar Transformation	21
5	Tuesday, September 20th 2025	22
5.1	Propagation of Errors for a Single Variable	22
5.2	Variance Propagation for a Single Variable	22
5.3	Propagation of Errors for Multiple Variables	23
5.4	Covariance Propagation for Functions of Multiple Variables	23
5.5	Examples of Error Propagation in Measurements	24
5.6	Matrix Formulation of Linear Error Propagation	24
5.7	Variance of the Arithmetic Mean	24
5.8	Example: Measuring the Period of a Sine Wave	25
6	Tuesday, September 25th 2025	26
6.1	Covariance Transformation Under Linear Transformations	26
6.2	The Binomial Distribution	27
7	Thursday, October 2nd 2025	30
7.1	Review of the Binomial Distribution and Its Properties	30
7.2	Bernoulli Distribution as a Special Case of the Binomial	30
7.3	Negative Binomial and Geometric Distributions	31
7.4	Samples and the Concept of an Ensemble	31
7.5	Poisson Distribution as a Limit of the Binomial	32
7.6	Poisson Process and Radioactive Decay	32
7.7	Moments and Variance of the Poisson Distribution	33
8	Tuesday, October 7th, 2025	34
8.1	The Gaussian (Normal) Distribution	34
8.2	The Standard Normal Distribution	34
8.3	Example: Circular Symmetry (Darts on a Board)	34
8.4	Expectation Value of a Gaussian	34
8.5	Moments of the Gaussian	35
8.6	Variance of the Gaussian	35
8.7	Kurtosis of the Gaussian	35
8.8	Poisson Distribution and Gaussian Limit	35
8.9	Central Limit Theorem (CLT)	36
8.10	Cumulative Distribution Function of a Gaussian	36
8.11	Gaussian Confidence Intervals	36
8.12	Estimators	36
9	Thursday, October 9th 2025	37
9.1	Properties of Estimators	37
9.2	Example: Measurements and Models	37
9.3	Least Squares Estimation	37
9.4	Straight Line Fit	37
9.5	Generalized Least Squares with Covariance Matrix	38
9.6	Covariance of the Estimated Parameters	38
9.7	Goodness of Fit	38

10 Tuesday, October 14th 2025	39
10.1 Chi-Squared for Uncorrelated and Correlated Measurements	39
10.2 Covariance Matrix and Linear Transformations	39
10.3 Modeling Data with Parameters	39
10.4 Least Squares Estimation	39
10.5 Distribution of Parameter Estimates	40
10.6 Quadratic Expansion of chi-squared and Error Estimates	40
10.7 Covariance of Parameter Estimates	40
11 Thursday, October 16th 2025	41
11.1 Chi-Squared Minimization and Degrees of Freedom	41
11.2 Two-Measurement Example and Correlated Variables	41
11.3 Residuals and Goodness of Fit	42
11.4 Distribution of Estimators	42
11.5 Toy Monte Carlo Simulations for Estimator Distributions	42
11.6 Typical Applications of Least Squares Fitting	43
12 Tuesday, October 21st 2025	45
12.1 Least-Squares Fits	45
12.2 Unbinned Data and Likelihood Functions	45
12.3 Example: Exponential Distribution	46
12.4 Example: Lifetime with Cutoff T	46
12.5 Example: Multiple Gaussian Measurements	47
12.6 Properties of the Maximum-Likelihood Estimator	47
13 Thursday, October 23rd 2025	48
13.1 Definition of the Likelihood Function	48
13.2 Maximum Likelihood Estimation (MLE)	48
13.3 Quadratic Approximation of the Log-Likelihood	48
13.4 Asymptotic Limit and Expectation Relation	48
13.5 Normalization of the Likelihood Function	49
13.6 Gaussian Approximation via the Central Limit Theorem	49
13.7 Variance of the Estimator and the Fisher Information	49
13.8 Taylor Expansion Near the Maximum Likelihood Estimate	50
13.9 Goodness of Fit and the Kolmogorov–Smirnov Test	51
14 Tuesday, October 28th 2025	52
14.1 Taylor Expansion of the Log-Likelihood Around the Maximum	52
14.2 Parameter of Interest and Nuisance Parameters	52
14.3 Likelihood versus Chi-Squared Interpretation	52
14.4 Normalization and Extended Likelihood	52
14.5 Nuisance Parameters and Systematic Effects	53
14.6 Gaussian Data Likelihood and Chi-Squared Minimization	53
14.7 Introduction to Hypothesis Testing	54
14.8 Rejection Regions and Type I Error	54
14.9 Neyman–Pearson Lemma and Type II Error	54
14.10 Likelihood Ratio Test and Interpretation	55

15 Thursday, October 30th 2025	56
15.1 Introduction to Hypothesis Testing	56
15.2 The Neyman-Pearson Lemma and Optimal Test Regions	57
15.3 Gaussian Example: Testing Mean Values	58
15.4 Applied Example: Mineral Density Classification	59
15.5 Common Misconceptions in Hypothesis Testing	59
16 Tuesday, November 4th 2025	60
16.1 Likelihood in Hypothesis Testing and Nuisance Parameters	60
16.2 Bayes' Theorem and Its Application to Hypothesis Testing	60
16.3 Maximum A Posteriori (MAP) Estimation	61
16.4 MAP Estimation for Normal Distribution with Gaussian Prior	61
16.5 Conjugate Priors and the Beta Distribution	62
16.6 Choice of Prior: Practical Considerations and Pitfalls	62
17 Thursday, November 6th 2025	63
17.1 Deductive and Inductive Reasoning; Weak Syllogism and Bayesian Motivation	63
17.2 Plausibility Notation and Consistency Requirements	64
17.3 Product Rule for Combined Plausibility	64
17.4 Solution of the Associativity Equation via Transform	65
18 Thursday, November 13th 2025	66
18.1 Plausibility as a Generalization of Logic	66
18.2 Derivation of the Functional Form via Associativity	66
18.3 Product Rule for Plausibility	68
18.4 Sum Rule for Plausibility	69
18.5 Identification of Plausibility with Probability	69
18.6 Application to Deductive Logic	70
18.7 Mutually Exclusive and Exhaustive Propositions	71
19 Tuesday, November 18th 2025	72
19.1 Continuous Variables: CDF and PDF	72
19.2 Priors and Conjugate Families	72
19.3 Fisher Information and Reparameterization	73
19.4 Jeffreys Prior: Bernoulli Trial Example	75
19.5 Maximum Entropy Principle	75
20 Thursday, November 20th 2025	77
20.1 Overview of Prior Distributions	77
20.2 Maximum Entropy Prior with Known Mean	77
20.3 Maximum Entropy for Continuous Variables	78
20.4 Scale Invariant and Location Invariant Priors	79
20.5 Nuisance Parameters and Marginalization	80
21 Tuesday November 25th 2025	81
21.1 Inverse Transform Sampling	81
21.2 Example: Sine Distribution	82
21.3 Rejection Sampling (Accept-Reject Method)	82
21.4 Basic Rejection Sampling with Uniform Proposal	83
21.5 General Rejection Sampling with Arbitrary Proposal	84

22 Thursday, November 27th 2024	85
22.1 Curse of Dimensionality	85
22.2 Monte Carlo Integration	85
22.3 Importance Sampling	86
22.4 Importance Sampling for Bayesian Inference	87
22.5 Metropolis-Hastings Algorithm	88

1 Tuesday, September 9th 2025

This lecture covered the course structure and grading.

1.1 Motivation: Stochastic Nature of Experimental Data

- Stochastic processes:
 - muon decay
 - inherent stochasticity
 - quantum mechanics
- Mostly concerned with measurement devices — how do we measure?
- Example: a muon lifetime experiment
 - Take a cosmic muon, detect light, and discriminate.
 - Muon decays into an electron and neutrinos, and the electron produces light.
 - Measure the time between light pulses.
 - Many factors cause noise in the data — results change even if the same mechanism occurs twice.

1.2 Probabilistic Interpretation of Experimental Results

- Experiments are repeated trials.
- Probability (probabilistic interpretation):
 - Results are interpreted as the long-term average of repeating an experiment many times.
 - Example: coin flip

$$P(H) = \lim_{N \rightarrow \infty} \frac{n_H}{N}$$

$n(H)$ = number of heads in N trials

1.3 Sample Spaces and Stochastic Variables

- In modern probability theory:
 - 3 axioms (Kolmogorov)
 - Let X be a stochastic variable.
 - Define sample space S (Ω):

$$S = \{x_1, x_2, \dots\}$$

- Examples:

1. Coin flip:

$$S = \{H, T\}$$

2. Roll a die:

$$S = \{1, 2, 3, 4, 5, 6\}$$

3. Grade in this class:

$$S = \{0, 1, 2, \dots, 100\}$$

4. Decay time of a radioactive atom:

$$S = [0, \infty)$$

- S can be finite (Binomial), countable (Poisson), or infinite (Gaussian, Uniform).

1.4 Events and Set Operations

- Definition: An event E is a subset of S .
- Example: one die roll

$$S = \{1, 2, 3, 4, 5, 6\}$$

$$E = \text{rolling an even number} = \{2, 4, 6\}$$

- Example: E = atom decayed by time t_0

$$S = [0, t_0]$$

- Operations on events:
 - Union (OR) and Intersection (AND)
 - Let A, B be events in S :

$$E = A \cup B = \{e : e \in A \text{ or } e \in B \text{ (or both)}\}$$

- Example: flip a coin twice

$$S = \{HH, HT, TH, TT\}$$

$$A = \text{1st flip is H} = \{HH, HT\}$$

$$B = \text{2nd flip is H} = \{HH, TH\}$$

$$A \cup B = \{HH, HT, TH\}$$

$$A \cap B = \{e \mid e \in A \text{ and } e \in B\} = \{HH\}$$

$$AB = A \cap B$$

$$A^c = \{e \mid e \in S \text{ and } e \notin A\} = \{TH, TT\}$$

- Properties:

- Commutative:

$$A \cup B = B \cup A, \quad AB = BA$$

- Associative:

$$A \cup (B \cup C) = (A \cup B) \cup C, \quad (AB)C = A(BC)$$

- Distributive:

$$(A \cup B)C = AC \cup BC, \quad A(B \cup C) = AB \cup AC$$

- De Morgan's Laws:

$$(A \cup B)^c = A^c B^c, \quad (AB)^c = A^c \cup B^c$$

1.5 Kolmogorov's Axioms of Probability

- A function P on S is a probability measure if it satisfies:

1. $P(S) = 1$
2. $P(\emptyset) = 0$
3. For any countable sequence of disjoint events E_1, E_2, \dots in S :

$$E_i E_j = \emptyset \text{ for } i \neq j$$

$$P\left(\bigcup_{i=1}^{\infty} E_i\right) = \sum_{i=1}^{\infty} P(E_i)$$

1.6 Consequences of the Probability Axioms

•

$$P(\emptyset) = 0$$

Let

$$E_1 = S, \quad E_2 = \emptyset$$

$$E_1 E_2 = \emptyset$$

$$P(S \cup \emptyset) = P(S) + P(\emptyset) = 1 + P(\emptyset)$$

$$P(S) = 1, P(\emptyset) = 0$$

•

$$P(E^c) = 1 - P(E)$$

$$1 = P(S) = P(E \cup E^c) = P(E) + P(E^c)$$

• If $B \subset A$, then:

$$P(B) \leq P(A)$$

$$A = B \cup (B^c A)$$

$$P(A) = P(B \cup (B^c A))$$

$$P(B) = P(A) - P(B^c A) \leq P(A)$$

•

$$P(A \cup B) = P(A) + P(B) - P(AB)$$

If we let the areas of the Venn diagram be 1 (A), 2 (A+B), 3 (B), then:

$$A \cup B = 1 \cup 2 \cup 3$$

$$P(A \cup B) = P(1 \cup 2 \cup 3) = P(1) + P(2) + P(3)$$

$$P(A) = P(1) + P(2), \quad P(B) = P(2) + P(3)$$

$$P(A) + P(B) - P(2) = P(1) + P(2) + P(3) = P(A \cup B)$$

$$\text{equivalently } P(A) + P(B) - P(AB) = P(A) + P(B) - P(AB)$$

1.7 Uniform Probability on Finite Sample Spaces

•

$$E_i = S_i \text{ for } i = 1, 2, \dots, n$$

$$E_i E_j = \emptyset \text{ for } i \neq j$$

$$S = \bigcup_{i=1}^n E_i$$

$$P(S) = 1 = P(\bigcup_{i=1}^n E_i) = \sum_{i=1}^n P(E_i)$$

$$P(E_i) = P(E_j) \quad \text{all equally likely}$$

$$1 = \sum_{i=1}^N P(E_i) = NP(E_i)$$

$$P(E_i) = \frac{1}{N} = P(E_j)$$

$N = |S|$ = number of elements in (cardinality of) S

F be any event (set) in S with k elements $|F| = k$

$$P(F) = P(\bigcup_{S_i \in F} \{E_i\}) = \sum_{i=1}^k P(E_j) = \sum_{i=1}^k \frac{1}{N} = \frac{k}{N} = \frac{|F|}{|S|}$$

1.8 Example: Probability of a Straight in Poker

- Example: 5-card poker hand forming a straight

$$S = \{(AC, 2C, 3C, 4C, 5C), (2C, 3C, 4C, 5C, 6C), \dots\}$$

$$S = \binom{52}{5} = \frac{52!}{5!47!} = 2,598,960$$

- Event = straight = 5 consecutive cards, not of the same suit, any starting card.

$$10(4^5 - 4) = 10200$$

- Starting cards: Ace (A,2,3,4,5), 2 (2,3,4,5,6), ..., 10 (10,J,Q,K,A)
- Not all the same suit: $4^5 - 4$ (exclude all same suit)

$$P(\text{straight}) = \frac{10(4^5 - 4)}{\binom{52}{5}} = 0.00392465$$

1.9 Conditional Probability

- Given 2 events E, F , sample space S :

$$P(E) = \text{probability of a trial from } S \text{ in } E$$

$$P(F) = \text{probability of a trial from } S \text{ in } F$$

- Conditional probability of E given F has occurred:

$$P(E|F) = \text{probability of a trial from } S \text{ in } E, \text{ given the trial is in } F$$

- Note: $P(EF)$ is the probability of a trial from S in both E and F .
- Need to normalize by $P(F)$, so we define:

$$P(E|F) = \frac{P(EF)}{P(F)} \quad \text{if } P(F) > 0$$

$$P(EF) = P(E|F)P(F)$$

- Example: flip a coin 2 times

$$S = \{HH, HT, TH, TT\}$$

Conditional probability of $HH \equiv A$ given:

- First flip = $H \equiv B = \{HH, HT\}$
- Either flip is $H \equiv C = \{HH, HT, TH\}$

$$P(A|B) = \frac{P(AB)}{P(B)} = \frac{P(\{HH\})}{P(\{HH, HT\})} = \frac{\frac{1}{4}}{\frac{1}{2}} = \frac{1}{2}$$

$$P(A|C) = \frac{P(AC)}{P(C)} = \frac{P(\{HH\})}{P(\{HH, HT, TH\})} = \frac{\frac{1}{4}}{\frac{3}{4}} = \frac{1}{3}$$

2 Thursday, September 11th 2025

2.1 Bayes' Formula

- Let E, F be events:

$$E = EF \cup EF^c$$

$$P(E) = P(EF) + P(EF^c)$$

$$P(E) = P(E|F)P(F) + P(E|F^c)P(F^c)$$

$$P(E) = P(E|F)P(F) + P(E|F^c)(1 - P(F))$$

- Example:** Suppose a blood test is 95% effective in detecting a disease if the person has it. It also has a 1% false positive rate. Suppose 0.5% of the population has the disease.

D = person has disease

E = test is positive

- We want:

$$P(D|E) = \frac{P(ED)}{P(E)}$$

$$\begin{aligned} P(D|E) &= \frac{P(E|D)P(D)}{P(E|D)P(D) + P(E|D^c)(1 - P(D))} \\ &= \frac{0.95 \times 0.005}{0.95 \times 0.005 + 0.01 \times 0.995} = 0.32 \end{aligned}$$

- So even with a positive test, there is only a 32% chance of having the disease.

2.2 Law of Total Probability

- Let $\{F_i\}$ be mutually exclusive events such that:

$$\cup_{i=1}^n F_i = S$$

Then for any event E :

$$E = E \cap (\cup_{i=1}^n F_i) = \cup_{i=1}^n (EF_i)$$

$$P(E) = P(\cup EF_i) = \sum_{i=1}^n P(EF_i) = \sum_{i=1}^n P(E|F_i)P(F_i)$$

2.3 Independent Events

- Generally, $P(E|F) \neq P(E)$.
- If knowing F does not change the probability of E :

$$P(E|F) = \frac{P(EF)}{P(F)} = P(E)$$

$$\boxed{P(EF) = P(E)P(F)}$$

2.3.1 Example: Rolling Two Dice

- Let:

$$E_1 \equiv \text{sum} = 6$$

$$F \equiv \text{first die} = 4$$

$$E_1 : \{(1,5), (2,4), (3,3), (4,2), (5,1)\}$$

$$F : \{(4,1), (4,2), (4,3), (4,4), (4,5), (4,6)\}$$

$$E_1 F = \{(4,2)\}$$

$$P(E_1 F) = \frac{1}{36}$$

$$P(E_1) = \frac{5}{36}$$

$$P(F) = \frac{6}{36} = \frac{1}{6}$$

$$P(E_1)P(F) = \frac{5}{36} \times \frac{1}{6} = \frac{5}{216} \neq P(E_1 F)$$

- Let:

$$E_2 \equiv \text{sum} = 7$$

$$E_2 : \{(1,6), (2,5), (3,4), (4,3), (5,2), (6,1)\}$$

$$E_2 F = \{(4,3)\}$$

$$P(E_2) = \frac{6}{36} = \frac{1}{6}$$

$$P(F) = \frac{1}{6}$$

$$P(E_2 F) = \frac{1}{36}$$

2.4 Random Variables and Probability Distributions

- $S = \{\text{all possible outcomes of stochastic process } X\}$

$x = \text{random variable}$

$S = \text{finite or countable infinite: discrete random variable}$

$S = \text{uncountable infinite: continuous random variable}$

- Continuous case:

$$P(x_0, x_0 + dx) = p(x)dx$$

where $p(x)$ is the probability density function (pdf).

- Discrete case:

$$S = S_i$$

$p_i = \text{probability of } S_i \quad (\text{probability mass function, pmf})$

$$0 \leq P(S_i) \leq 1$$

$$1 = P(S)$$

$$0 \leq p(x)$$

$$\int_{-\infty}^{\infty} p(x)dx = 1$$

2.5 Describing a Distribution

- To describe $p(x)$ in general we specify:
 - **Mode** — peak value of $p(x)$
 - **Median** — 50% cumulative value
 - **Mean** — average value of x weighted by $p(x)$

2.6 Cumulative Distribution Function (CDF)

•

$$F(x) = \int_{-\infty}^x p(x') dx' = P(X \leq x)$$

$$F(-\infty) = 0, \quad F(\infty) = 1$$

2.7 Expectation Values

- Expectation of any function $f(x)$ over $p(x)$:

$$E(f) = \int_{\Omega} f(x) p(x) dx$$

$$E \text{ is a linear operator: } E(af + bg) = aE(f) + bE(g)$$

- Expectation of powers of x :

$$E(x^0) = E(1) = \int 1 \cdot p(x) dx = 1$$

$$E(x^1) = \int x p(x) dx \equiv \mu = \text{mean value of } x$$

$$E(x^2) = \int x^2 p(x) dx \equiv \sigma^2 = \text{variance of } x$$

2.8 Characteristic Function

- The characteristic function of $p(x)$:

$$\varphi(t) = \int_{-\infty}^{\infty} e^{itx} p(x) dx = E(e^{itx})$$

$$\varphi(t) = E \left(1 + itx + \frac{(itx)^2}{2!} + \dots \right)$$

$$= 1 + itE(x) + \frac{(it)^2}{2!} E(x^2) + \dots$$

$$= \sum_{k=0}^{\infty} \frac{(it)^k}{k!} \mu_{k'}$$

- Moments from $\varphi(t)$:

$$\left. \frac{d^n \varphi(t)}{dt^n} \right|_{t=0} = i^n \mu_{n'}$$

2.9 Central Moments

-

$$E((x - \mu)^n) = \int (x - \mu)^n p(x) dx \equiv \mu_n$$
$$\mu = E(x)$$

- 1st central moment:

$$E((x - \mu)^1) = E(x) - E(\mu) = \mu - \mu = 0$$

- 2nd central moment (variance):

$$E((x - \mu)^2) \equiv V(x) = \sigma^2$$

- 3rd central moment (skewness):

$$\text{skewness} = \frac{E((x - \mu)^3)}{\sigma^3}$$

- 4th central moment (kurtosis):

$$\text{kurtosis} = \frac{E((x - \mu)^4)}{\sigma^4} - 3$$

(The -3 ensures that the kurtosis of a normal distribution is 0.)

3 Tuesday, September 16th 2025

3.1 Bayes Theorem and Its Applications

- Bayes Theorem: for events A and B , we have

$$P(AB) = P(A|B)P(B) = P(B|A)P(A) = P(BA)$$

- Usually it is given in this form:

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

- People argued about when you are allowed to use this theorem.

3.2 The Monty Hall Problem: A Bayesian Analysis

- Example: Monty Hall Problem (Game show with host named Monty Hall)
 - There are 3 doors; behind one is a car, behind the other two are goats.
 - You select a door; if the car is behind it, you win.
 - Twist: after you select a door, Monty opens one of the other 2 doors to reveal a goat.
 - Question: stay or switch?
 - Solution: use Bayes theorem.
 - Sample space: $S = \{C_1 = \text{cgg}, C_2 = \text{gcg}, C_3 = \text{ggc}\}$
 - Event 2 = MH opens door 2.
 - Event 3 = MH opens door 3.
 - Number such that your choice is door 1.
 - Take case E_2 , then we want to know $P(C_1|E_2)$.

$$P(C_1|E_2) = \frac{P(E_2|C_1)P(C_1)}{P(E_2)}$$

- $P(C_1) = \frac{1}{3}$
- $P(E_2|C_1) = \frac{1}{2}$ because if the car is behind door 1, Monty can open either door 2 or 3.
- $P(E_2) = \frac{1}{2}$
- Law of total probability:

$$P(E_2) = P(E_2|C_1)P(C_1) + P(E_2|C_2)P(C_2) + P(E_2|C_3)P(C_3) = \frac{1}{2} \cdot \frac{1}{3} + 0 \cdot \frac{1}{3} + 1 \cdot \frac{1}{3} = \frac{1}{2}$$

- $P(C_1|E_2) = \frac{P(E_2|C_1)P(C_1)}{P(E_2)} = \frac{\frac{1}{2} \cdot \frac{1}{3}}{\frac{1}{2}} = \frac{1}{3}$
- $P(C_1|E_2) = \frac{1}{3}$
- $P(C_2|E_2) = 0$
- $P(C_3|E_2) = \frac{P(E_2|C_3)P(C_3)}{P(E_2)} = \frac{1 \cdot \frac{1}{3}}{\frac{1}{2}} = \frac{2}{3}$

3.3 Alternate Monty Hall Formulations

- Alternate version: E = MH shows you a goat from $\{2, 3\}$.

- We want to find $P(C_1|E)$.

- $P(C_1|E) = \frac{P(E|C_1)P(C_1)}{P(E)}$

- $P(C_1) = \frac{1}{3}$

- $P(E|C_1) = 1$ because if the car is behind door 1, Monty can open either door 2 or 3.

- $P(E) = 1$ by law of total probability:

$$P(E) = P(E|C_1)P(C_1) + P(E|C_2)P(C_2) + P(E|C_3)P(C_3) = 1 \cdot \frac{1}{3} + 1 \cdot \frac{1}{3} + 1 \cdot \frac{1}{3} = 1$$

- $P(C_1|E) = \frac{1 \cdot \frac{1}{3}}{1} = \frac{1}{3}$

- Another version: What if MH does *not* know where the car is?

- E = MH opens $\{2, 3\}$ and reveals a goat.

- We want to find $P(C_1|E)$.

- $P(C_1|E) = \frac{P(E|C_1)P(C_1)}{P(E)}$

- $P(C_1) = \frac{1}{3}$ because we picked door 1.

- $P(E|C_1) = \frac{1}{2}$ because if the car is behind door 1, Monty can open either door 2 or 3 since he does not know where the car is.

- By law of total probability:

$$P(E) = P(E|C_1)P(C_1) + P(E|C_2)P(C_2) + P(E|C_3)P(C_3) = 1 \cdot \frac{1}{3} + \frac{1}{2} \cdot \frac{1}{3} + \frac{1}{2} \cdot \frac{1}{3} = \frac{2}{3}$$

- $P(C_1|E) = \frac{\frac{1}{3}}{\frac{2}{3}} = \frac{1}{2}$

3.4 Monty Hall Generalized to n Doors

- Now back to the standard version of the problem but with n doors.

- You pick door 1, MH opens any door with a goat behind it from 2 to n ($n - 1$ options).

- $P(C_1|E) = \frac{P(E|C_1)P(C_1)}{P(E)}$

- $P(E) = 1$ because he can always choose a door with a goat behind it (many options and he knows the answers).

- $P(C_1) = \frac{1}{n}$

- $P(E|C_1) = 1$ because if the car is behind door 1, Monty can open any of the other doors.

3.5 Continuous Probability Distributions and Moments

- Continuous probability distribution $p(x)$:
- Moments:

$$E(x^n) = \int_{-\infty}^{\infty} x^n p(x) dx$$

$$\begin{aligned} \text{mean: } & \mu = E(x) \\ \text{variance: } & V(x) = \sigma^2 = E((x - \mu)^2) = E(x^2) - \mu^2 \\ \text{std dev: } & \sigma = \sqrt{\sigma^2} \end{aligned}$$

- Central moments:

$$\begin{aligned} E(x - \mu) &= E(x) - \mu = 0 \\ E((x - \mu)^2) &= \sigma^2 \\ E((x - \mu)^3) &= \text{skewness} \\ E((x - \mu)^4) &= \text{kurtosis} \end{aligned}$$

- Characteristic function:

$$\Phi(t) = E(e^{itx}) = \int_{-\infty}^{\infty} e^{itx} p(x) dx \quad (3.1)$$

$$= \sum_{k=0}^{\infty} \frac{(it)^k}{k!} \mu_k \quad (3.2)$$

$$\Phi_{\mu}(t) = E(e^{it(x-\mu)}) = E(e^{itx})e^{-it\mu} = \Phi(t)e^{-it\mu}$$

$$V(x) = E((x - \mu)^2) \quad (3.3)$$

$$= E(x^2 - 2\mu x + \mu^2) \quad (3.4)$$

$$= E(x^2) - 2\mu E(x) + \mu^2 E(1) \quad (3.5)$$

$$= E(x^2) - 2\mu^2 + \mu^2 \quad (3.6)$$

$$= E(x^2) - \mu^2 = E(x^2) - (E(x))^2 \quad (3.7)$$

3.6 Discrete Probability Distributions

- The discrete case (e.g., rolling a die, picking a card) uses a probability mass function.
- Usually denote outcomes as r :
- p_r = probability of outcome r .
- $\sum_r p_r = 1$
- $E(r) = \sum_r p_r r$ = mean μ
- Variance: $V(r) = \sum_r (r - \mu)^2 p_r = E(r^2) - \mu^2$
- Coin flip example: $S = \{H, T\}$.
- Often map to 0 or 1: $H = 0, T = 1$.
- But in theory you can pick any two numbers a and b to map outcomes, just so you can calculate mean and variance.

$$E(r) = ap_H + bp_T$$

3.7 Cumulative Distribution Functions

- For continuous case:

$$F(x) = \int_{-\infty}^x f(x') dx'$$

- For discrete case:

$$F(r) = \sum_{r' \leq r} p_{r'}$$

- $F(x)$ is the cumulative distribution function (CDF).
- $F(x)$ is non-decreasing, $F(-\infty) = 0$, $F(\infty) = 1$.

3.8 Multivariate Distributions and Covariance

- Distribution of multiple variables:
- Elements belong to real vector space \mathbb{R}^n .
- $P(AB) \dots P(A, B)$
- $p(x_1, x_2, \dots, x_n) \geq 0$ is the joint probability distribution function (PDF).
- $\int_{\Omega} p(\vec{x}) d^n x = 1$
- $E(f(\vec{x})) = \int_{\Omega} f(\vec{x}) p(\vec{x}) d^n x$
- $\mu_i = \int x_i p(\vec{x}) d^n x$
- $V(x_i) = \sigma_i^2 = \int (x_i - \mu_i)^2 p(\vec{x}) d^n x$
- Covariance:
- $V_{i,j} = E((x_i - \mu_i)(x_j - \mu_j))$
- $V_{i,i} = \sigma_i^2 = E((x_i - \mu_i)^2)$ (variance)
- $V_{i,j} = V_{j,i}$ (symmetry)

4 Thursday, September 18th 2025

4.1 Conditional Probability: A Simple Example

- For fun, example that depends on cultural assumptions: A king comes from a family with two kids. What is the probability that the king's sibling is a sister?
- $S = \{(m, m), (m, f), (f, m), (f, f)\}$
- $P(S|K) = \frac{P(SK)}{P(K)} = \frac{\frac{1}{2}}{\frac{3}{4}} = \frac{2}{3}$

4.2 Distributions of Multiple Random Variables

- $p(x_1, x_2, \dots, x_n)$
- $S = \mathbb{R}^n$
- $\int p(x_1, x_2, \dots, x_n) dx_1 dx_2 \dots dx_n = 1$
- For any function $f(\vec{x})$:

$$E(f) = \int f(\vec{x})p(\vec{x})d\vec{x}$$

- $E(x_1) = \int x_1 p(\vec{x})d\vec{x} = \mu_1$
- $E(x_i) = \mu_i$
- $V(x_i) \equiv \sigma_i^2 = \int (x_i - \mu_i)^2 p(\vec{x})d\vec{x}$

4.3 Covariance Matrix and Correlation Coefficient

- Define covariance:

$$V_{ij} = E((x_i - \mu_i)(x_j - \mu_j))$$

- $V_{ii} = \sigma_i^2$ (variance)
- $V_{ij} = V_{ji}$ (symmetry)
- $V_{ij} = 0$ for independent variables
- Expanding the covariance matrix:

$$\begin{aligned} V_{ij}(\vec{x}) &= E((x_i - \mu_i)(x_j - \mu_j)) \\ &= E(x_i x_j - \mu_i x_j - \mu_j x_i + \mu_i \mu_j) \\ &= E(x_i x_j) - \mu_i E(x_j) - \mu_j E(x_i) + \mu_i \mu_j \\ &= E(x_i x_j) - \mu_i \mu_j - \mu_j \mu_i + \mu_i \mu_j \\ &= E(x_i x_j) - \mu_i \mu_j \end{aligned}$$

- So we can say that $V_{ij} \geq 0$
- V_{ij} can be negative, zero, or positive
- Define the correlation coefficient:

$$\rho(x_i, x_j) = \rho_{ij} = \frac{V_{ij}}{\sqrt{V_{ii}}\sqrt{V_{jj}}} = \frac{V_{ij}}{\sigma_i \sigma_j}$$

- We find that $-1 \leq \rho_{ij} \leq 1$

4.4 Independence and Uncorrelated Variables

- Random variables x_1, \dots, x_n are independent if the joint pdf factorizes:

$$p(x_1, \dots, x_n) = p_1(x_1)p_2(x_2) \dots p_n(x_n)$$

- Independent variables are uncorrelated:

$$\begin{aligned} E(x_i x_j) &= \int x_i x_j p(\vec{x}) d\vec{x} \\ &= \int x_i x_j p_1(x_1) \dots p_n(x_n) dx_1 \dots dx_n \\ &= \int x_i p_i(x_i) dx_i \int x_j p_j(x_j) dx_j \int p_2(x_2) dx_2 \dots \int p_n(x_n) dx_n = \mu_i \mu_j \end{aligned}$$

$$V_{ij} = E(x_i x_j) - \mu_i \mu_j$$

In the case of independent variables:

$$V_{ij} = \mu_i \mu_j - \mu_i \mu_j = 0$$

- Independent variables are uncorrelated, but uncorrelated variables are not necessarily independent.

4.5 Examples of Correlated and Uncorrelated Variables

- 100% correlation example:
- $x = \text{Uniform}[-1, 1]$, plot distribution from -1 to 1 .
- $y = x$:
- $V_{ij} = E(xy) - E(x)E(y) = E(x^2) = \int_{-1}^1 x^2 \frac{1}{2} dx = \frac{1}{3} \neq 0$
- $y = |x|$:
- $E(xy) = \int_{-1}^0 x(-x)p(x)dx + \int_0^1 x p(x)dx$
- $E(xy) = \int_0^1 x^2 \frac{1}{2} dx - \int_{-1}^0 x^2 \frac{1}{2} dx = \frac{1}{6} - \frac{1}{6} = 0$

4.6 Marginal Distributions

- For a joint pdf $p(x_1, x_2, \dots, x_n)$, the marginal probability density functions are:

$$f_1(x_1) = \int p(x_1, x_2, \dots, x_n) dx_2 dx_3 \dots dx_n$$

- If variables are independent:

$$\begin{aligned} f_1(x_1) &= \int p(x_1, x_2, \dots, x_n) dx_2 \dots dx_n \\ &= p_1(x_1) \int p_2(x_2) dx_2 \int p_3(x_3) dx_3 \dots \int p_n(x_n) dx_n \\ &= p_1(x_1) \cdot 1 \cdot 1 \cdot \dots \cdot 1 = p_1(x_1) \end{aligned}$$

4.7 Change of Variables in Probability Densities

- Something we need to know, because we do it all the time:
 - Change of variables of P
 - Calculate new V_{ij} under new variables
- Let x be a random variable with pdf $f(x)$ and let y be some function.
- First: y is one-to-one with f

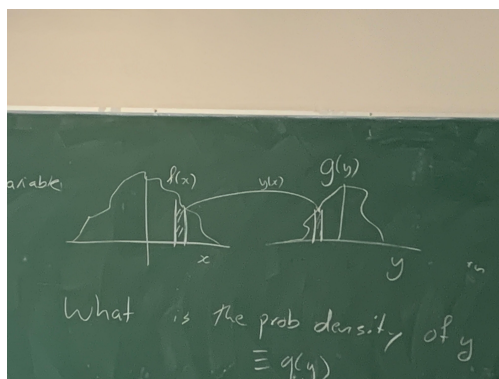


Figure 4.1: 1-to-1 function

- What is the probability density of y , denoted $g(y)$?
- Conservation of probability:
- $f(x)dx = g(y)dy$
- $g(y) = f(x) \left| \frac{dx}{dy} \right|$

$$f(x) \left| \frac{dx}{dy} \right| = g(y)$$

4.8 Change of Variables: Non One-to-One Case

- If y is not one-to-one: sum over all segments that map to the same y .
- Example: $f(x)$ uniform on $[0, 1]$, $f(x) = 1$
- Let $y(x) = \frac{-1}{\lambda} \ln(x)$
- $\frac{dy}{dx} = \frac{-1}{\lambda x}$
- $\frac{dx}{dy} = -\lambda x$
- $-\lambda x = \ln x$
- $e^{-\lambda y} = x$
- $\lambda > 0 \Rightarrow \frac{dx}{dy} = -\lambda x = -\lambda e^{-\lambda y}$
- $g(y) = f(x) \left| \frac{dx}{dy} \right| = 1 \cdot \lambda e^{-\lambda y} = \lambda e^{-\lambda y}$

4.9 Multivariate Transformations and the Jacobian

- If we have variables $\{x_i\}$ and transform to new variables $\{y_i\}$:
- Region \mathbb{R} in x -space maps to region \mathbb{R}' in y -space.

$$\int_{\mathbb{R}} f(\vec{x}) d\vec{x} = \int_{\mathbb{R}'} f(\vec{x}(\vec{y}))(\vec{y}) \left| \frac{\partial \vec{x}}{\partial \vec{y}} \right| d\vec{y}$$

$$g(\vec{y}) = f(\vec{x}(\vec{y})) |J|$$

- Where $\left| \frac{\partial \vec{x}}{\partial \vec{y}} \right|$ is the Jacobian determinant of the transformation.
- Jacobian matrix J :

$$J = \begin{bmatrix} \frac{\partial x_1}{\partial y_1} & \frac{\partial x_1}{\partial y_2} & \cdots & \frac{\partial x_1}{\partial y_n} \\ \frac{\partial x_2}{\partial y_1} & \frac{\partial x_2}{\partial y_2} & \cdots & \frac{\partial x_2}{\partial y_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial x_n}{\partial y_1} & \frac{\partial x_n}{\partial y_2} & \cdots & \frac{\partial x_n}{\partial y_n} \end{bmatrix}$$

4.10 Example: Cartesian to Polar Transformation

- Change to polar coordinates:
- $x = r \cos \theta$
- $y = r \sin \theta$
- $P'(r, \theta) = ? = p(x, y) \left| \frac{\partial(x, y)}{\partial(r, \theta)} \right|$
- $\frac{\partial x}{\partial r} = \cos \theta$
- $\frac{\partial y}{\partial r} = \sin \theta$
- $\frac{\partial x}{\partial \theta} = -r \sin \theta$
- $\frac{\partial y}{\partial \theta} = r \cos \theta$
- $J = \begin{bmatrix} \cos \theta & -r \sin \theta \\ \sin \theta & r \cos \theta \end{bmatrix}$
- $J = r \cos^2 \theta + r \sin^2 \theta = r$
- $p'(r, \theta) = \frac{r}{\pi} dr d\theta$

5 Tuesday, September 20th 2025

5.1 Propagation of Errors for a Single Variable

- Given $f(x)$ pdf, $\mu \equiv E(x)$, $\sigma^2 \equiv V(x) = E(x^2) - \mu^2$
- Know $f(x) \rightarrow g(y)$, given $y(x)$.
- Taylor expand $y(x)$ about mean μ :

$$y(x) = y(\mu) + y'(\mu)(x - \mu) + \frac{1}{2!}y''(\mu)(x - \mu)^2 + \dots$$

$$E(y(x)) \equiv \mu_y$$

$$\begin{aligned} E(y(x)) &= E(y(\mu)) + y'(\mu)E(x - \mu) + \frac{1}{2!}y''(\mu)E((x - \mu)^2) + \dots \\ &= y(\mu) + y'(\mu) \cdot 0 + \frac{1}{2!}y''(\mu)V(x) + \dots \end{aligned}$$

- To the 1st order:

$$\mu_y = E(y(x)) = y(\mu) = y(E(x))$$

5.2 Variance Propagation for a Single Variable

- Variance of y :

$$V(y) = E((y(x) - E(y(x)))^2) \quad (5.1)$$

$$= E((y(x) - \mu_y)^2) \quad (5.2)$$

$$= E((y'(\mu)(x - \mu) + \frac{1}{2!}y''(\mu)(x - \mu)^2 + \dots)^2) \quad (5.3)$$

$$= E(y'(\mu)^2(x - \mu)^2 + y'(\mu)y''(\mu)(x - \mu)^3 + O((x - \mu)^4)) \quad (5.4)$$

$$= y'(\mu)^2V(x) + \dots \quad (5.5)$$

- Some relations:

$$E(x) \equiv \mu_x$$

$$V(x) \equiv \sigma_x^2$$

$$y = y(x)$$

$$E(y) \equiv \mu_y = y(\mu_x)$$

$$V(y) \equiv \sigma_y^2 = (y'(\mu_x))^2\sigma_x^2$$

$$\sigma_y = |y'(\mu_x)|\sigma_x$$

- Example: $y = \frac{1}{x}$, $\frac{dy}{dx} = -\frac{1}{x^2}$

$$\sigma_y^2 = \frac{1}{\mu_x^4}\sigma_x^2$$

5.3 Propagation of Errors for Multiple Variables

- Let us suppose we have n variables $\{x_i\}$, with pdf $f(\vec{x})$.
- Let $y_j = 1, 2, \dots, m$ be m functions of x_i .
- $y_j = y_j(x_1, x_2, \dots, x_n)$
- $V_{ij}(x)_{n \times n}(\vec{x}) = \text{covariance matrix of } \{x_i\}$
- $V_{ij}(\vec{x}) = E((x_i - \mu_i)(x_j - \mu_j))$
- Taylor expand each y_j : $\vec{\mu} = (\mu_1, \mu_2, \dots, \mu_n)$
- $y_j(\vec{x}) = y_j(\vec{\mu}) + \sum_i \frac{\partial y_j}{\partial x_i} \Big|_{\vec{\mu}} (x_i - \mu_i) + \frac{1}{2!} \sum_{i,k} \frac{\partial^2 y_j}{\partial x_i \partial x_k} \Big|_{\vec{\mu}} (x_i - \mu_i)(x_k - \mu_k) + \dots$
- $E(y_j(\vec{x})) = E(y_j(\vec{\mu})) + \sum \frac{\partial y_j}{\partial x_i} E(x_i - \mu_i) + \dots = y_j(\vec{\mu})$

5.4 Covariance Propagation for Functions of Multiple Variables

- Covariance between y_k and y_l :

$$\begin{aligned}
 & E((y_k - \mu_{y_k})(y_l - \mu_{y_l})) \\
 &= E((y_k - y_k(\mu))(y_l - y_l(\mu))) \\
 &= E\left(\sum_i \frac{\partial y_k}{\partial x_i} \Big|_{\mu} (x_i - \mu_i) \sum_j \frac{\partial y_l}{\partial x_j} \Big|_{\mu} (x_j - \mu_j)\right) \\
 &= \sum_{i,j} \frac{\partial y_k}{\partial x_i} \Big|_{\mu} \frac{\partial y_l}{\partial x_j} \Big|_{\mu} E((x_i - \mu_i)(x_j - \mu_j))
 \end{aligned}$$

$$\boxed{V_{kl}(\vec{y})_{m \times m} = \sum_{i,j} \frac{\partial y_k}{\partial x_i} \Big|_{\vec{\mu}} \frac{\partial y_l}{\partial x_j} \Big|_{\vec{\mu}} V_{ij}(\vec{x})_{n \times n}}$$

- Example: x, y random variables,

$$V(x, y) = \begin{bmatrix} V_{xx} & V_{xy} \\ V_{yx} & V_{yy} \end{bmatrix} = \begin{bmatrix} \sigma_x^2 & \rho_{xy} \sigma_x \sigma_y \\ \rho_{xy} \sigma_x \sigma_y & \sigma_y^2 \end{bmatrix}$$

- $z = x + y$
- $V(z) = \sigma_z^2 = \left(\frac{\partial z}{\partial x}\right)^2 V_{xx} + 2 \frac{\partial z}{\partial x} \frac{\partial z}{\partial y} V_{xy} + \left(\frac{\partial z}{\partial y}\right)^2 V_{yy}$
- $= \sigma_x^2 + 2\rho_{xy} \sigma_x \sigma_y + \sigma_y^2$
- If x_i are uncorrelated,

$$\begin{aligned}
 V_{ij} &= \sigma_{i,j} \sigma_i^2 = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix} \\
 V_{kl}(\vec{y}) &= \sum_i \frac{\partial y_k}{\partial x_i} \Big|_{\mu} \frac{\partial y_l}{\partial x_i} \Big|_{\mu} V_{ii}(\vec{x}) \\
 \text{variance } V_{kk} &= \sum_i \left(\frac{\partial y_k}{\partial x_i}\right)^2 \sigma_i^2
 \end{aligned}$$

5.5 Examples of Error Propagation in Measurements

- Example: Measuring resistances. x_i independent, $z = x + y$, $x = R_1$ resistor value, $y = R_2$ resistor value, $z = R_{\text{tot}}$ total resistance.
- $R_1 \pm \sigma_{R_1}$
- Convention is to use $\sqrt{V(R)}$ as uncertainty.
- For a good measuring device, $E(R) = R_{\text{true}} \leftarrow$ unbiased.
- $V(R) = \text{small}$
- $R_1 \pm \sigma_{R_1}$, $R_2 \pm \sigma_{R_2}$, then $\sigma_{R_{\text{tot}}} = \sqrt{\sigma_{R_1}^2 + \sigma_{R_2}^2}$
- $R = R_{\text{tot}} = R_1 + R_2$
- $z = xy$, like I, R
- $\sigma_z^2 = \left(\frac{\partial z}{\partial x}\right)^2 \sigma_x^2 + \left(\frac{\partial z}{\partial y}\right)^2 \sigma_y^2 = y^2 \sigma_x^2 + x^2 \sigma_y^2$

$$\left(\frac{\sigma_z}{z}\right)^2 = \left(\frac{\sigma_x}{x}\right)^2 + \left(\frac{\sigma_y}{y}\right)^2$$

5.6 Matrix Formulation of Linear Error Propagation

- Formula is exact if transformation of variables is linear.
- $\vec{y} = A\vec{x}$, A is $m \times n$ matrix, \vec{x} is $n \times 1$, \vec{y} is $m \times 1$.
- $\frac{\partial y_k}{\partial x_i} = \text{constant} \Rightarrow$ higher order terms in Taylor expansion are 0
- $V_{kl}(\vec{y}) = \sum_{i,j} \frac{\partial y_k}{\partial x_i} \frac{\partial y_l}{\partial x_j} V_{ij}(\vec{x})$
- Matrix notation:
- $V_{kl}(\vec{y}) = \sum_{i,j} A_{ki} A_{lj} V_{ij}(\vec{x})$
- $= \sum_{i,j} A_{ki} V_{ij}(\vec{x}) A_{lj}$
- $= \sum_{i,j} A_{ki} V_{ij} (A^T)_{jl}$
- $= (AV(\vec{x})A^T)_{kl}$

$$V(\vec{y})_{m \times m} = A_{m \times n} V(\vec{x})_{n \times n} A_{n \times m}^T$$

5.7 Variance of the Arithmetic Mean

- Example: Arithmetic mean. Let $x_i = n$ identical independent variables with $V(x_i) = \sigma_x^2$
- Set $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$
- Recall that $V(ax) = a^2 V(x)$
- $V(\bar{x}) = V\left(\frac{1}{n} \sum_{i=1}^n x_i\right) = \frac{1}{n^2} V\left(\sum_{i=1}^n x_i\right) = \frac{1}{n^2} \sum_{i=1}^n V(x_i) = \frac{1}{n^2} n \sigma_x^2 = \frac{\sigma_x^2}{n}$

- If variables are different σ_i^2 : n measurements
- $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$
- $V(\bar{x}) = \frac{1}{n^2} \sum_{i=1}^n V(x_i) = \frac{1}{n^2} \sum_{i=1}^n \sigma_i^2$
- $\sigma_{\bar{x}} = \frac{1}{n} \sqrt{\sum_{i=1}^n \sigma_i^2}$

5.8 Example: Measuring the Period of a Sine Wave

- Example: Measure period of sine wave on scope.
- $T = \Delta t = t_2 - t_1$
- $\sigma_T^2 = \left(\frac{\partial \Delta t}{\partial t_1} \right)^2 \sigma_t^2 + \left(\frac{\partial \Delta t}{\partial t_2} \right)^2 \sigma_t^2 = \sigma_t^2 + \sigma_t^2 = 2\sigma_t^2$
- Measure N cycles, $T = \frac{1}{N} \Delta t$
- $\sigma_{T^2} = \frac{1}{N^2} \sigma_{\Delta t}^2 = \frac{2}{N^2} \sigma_t^2$

6 Tuesday, September 25th 2025

6.1 Covariance Transformation Under Linear Transformations

- Linear transformation:

$$\vec{y} = A\vec{x}$$

$$V_{kl}(\vec{y}) = \sum_{i,j} \frac{\partial y_k}{\partial x_i} \frac{\partial y_l}{\partial x_j} V_{ij}(\vec{x})$$

- Linear $y_k = \sum A_{kj}x_j$
- then

$$V_{kl}(\vec{y}) = \sum_{i,j} A_{ki} A_{lj} V_{ij}(\vec{x})$$

- or in matrix form

$$V(\vec{y}) = \left(AV(\vec{x})A^T \right)_{kl}$$

Diagonalization via Eigenvectors

- If \hat{e}_i are the eigenvectors of V , then

$$V(\vec{x})\hat{e}_i = \lambda_i \hat{e}_i$$

- Form:

$$A = \begin{pmatrix} \hat{e}_1 \\ \dots \\ \hat{e}_n \end{pmatrix} = \begin{pmatrix} \hat{e}_{11} & \hat{e}_{12} & \dots & \hat{e}_{1n} \\ \dots & \dots & \dots & \dots \\ \hat{e}_{n1} & \hat{e}_{n2} & \dots & \hat{e}_{nn} \end{pmatrix}$$

- then

$$A^T A = I$$

- then:

$$VA^T = V \begin{pmatrix} \hat{e}_1 & \dots \\ \dots & \dots \\ \hat{e}_n & \dots \end{pmatrix} = \begin{pmatrix} \lambda_1 \hat{e}_{11} & \dots & \lambda_n \hat{e}_{n1} \\ \dots & \dots & \dots \\ \lambda_1 \hat{e}_{1n} & \dots & \lambda_n \hat{e}_{nn} \end{pmatrix}$$

- Then:

$$AVA^T = \begin{pmatrix} \hat{e}_{11} & \dots & \hat{e}_{1n} \\ \dots & \dots & \dots \\ \dots & \dots & \dots \end{pmatrix} \begin{pmatrix} \lambda_1 \hat{e}_{11} & \dots & \dots \\ \dots & \dots & \dots \\ \lambda_1 \hat{e}_{1n} & \dots & \dots \end{pmatrix} = \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & \lambda_n \end{pmatrix}$$

- Then:

$$AVA^T = V(\vec{y}) = \begin{pmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & \sigma_n^2 \end{pmatrix}$$

6.2 The Binomial Distribution

- Consider an experiment with two outcomes.
- E.g. coin flips, selecting a ball with 2 possible colours, etc.
- One trial is called a Bernoulli trial.

Bernoulli Trials and Sampling Methods

- Example – Method 1: You have an urn filled with N balls. Some are red (R), some are blue (B).
- (0) What is your estimate of n_R , n_B , or $f = n_R/N$ or p of drawing R?
- (1) You pick a ball: R. Q: estimate of $p = n_R/N$?
- (2) You pick another without replacing 1st ball: get R.
- (3) R
- (4) Get B
- This is a question about this ONE urn.
- Now Method 2: you draw red, and you PUT IT BACK. You repeat this several times.
- Now Method 3: We have an infinite source of balls with fraction p red and $(1 - p)$ blue.

$$P(R) = p$$
$$P(B) = 1 - p$$

Derivation of the Binomial Probability

- Make infinite number of urns all with N balls, with fraction p red and $(1 - p)$ blue.
- Open all, count n_R red balls, n_B blue balls.
- In our case we have N balls, prob $p = R$ and $1 - p = q = B$.
- Prob of getting sequence RRB is:

$$P(RRB) = p \cdot p \cdot (1 - p) = p^2(1 - p)$$

- If we don't care about order, then:

$$P(RRB) = P(RBR) = P(BRR) = p^2(1 - p)$$

- There are 3 ways of ordering RRB, so total probability is:

$$P(2R, 1B) = 3p^2(1 - p) = 3p^2q$$

- Number of ways to choose r items from N is:

$$\binom{N}{r} = \frac{N!}{r!(N - r)!}$$

- Probability of getting exactly r R out of N :

$$P_r = \binom{N}{r} p^r (1 - p)^{N-r} = B(r; N, p)$$

- This is called the Binomial distribution and applies to anything where there are 2 outcomes (A, \bar{A}).

Mean and Variance of the Binomial Distribution

- Want mean, σ

$$\begin{aligned}
 E(r) &= \sum_{r=0}^n r P_r = \sum_{r=0}^n r \binom{n}{r} p^r (1-p)^{n-r} \\
 &= \sum_{r=0}^n r \frac{n!}{r!(n-r)!} p^r (1-p)^{n-r} \\
 &= \sum_{r=1}^n \frac{n!}{(r-1)!(n-r)!} p^r (1-p)^{n-r} \\
 &= np \sum_{r=1}^n \frac{(n-1)!}{(r-1)!(n-r)!} p^{r-1} (1-p)^{n-r}
 \end{aligned}$$

- Change sum $r' = r - 1 \rightarrow n' = n - 1$

$$\begin{aligned}
 E(r) &= np \sum_{r'=0}^{n-1} \frac{(n-1)!}{r'!(n-1-r')!} p^{r'} (1-p)^{(n-1)-r'} \\
 &= np \sum_{r'=0}^{n-1} \binom{n-1}{r'} p^{r'} (1-p)^{n-1-r'} \\
 &= np \cdot 1 = np
 \end{aligned}$$

- from:

$$\begin{aligned}
 (p+q)^n &= \sum_{r=0}^n \binom{n}{r} p^r q^{n-r} \\
 (p+1-q)^n &= 1^n = 1 \\
 E(r) &= np
 \end{aligned}$$

- This is what we want!
- Now:

$$V(r) = \sum r^2 p_r - E(r)^2 = \sum r^2 p_r - n^2 p^2$$

- Slightly easier to calculate:

$$\begin{aligned}
 \sum r(r-1) p_r &= \sum_{r=0}^n r(r-1) \frac{n!}{r!(n-r)!} p^r (1-p)^{n-r} = \sum_{r=2}^n \frac{n!}{(r-2)!(n-r)!} p^r q^{n-r} \\
 &= n(n-1) p^2 \sum_{r=2}^n \frac{(n-2)!}{(r-2)!(n-r)!} p^{r-2} q^{n-r}
 \end{aligned}$$

- Sub $r' = r - 2$

$$\begin{aligned}
 &= n(n-1) p^2 \sum_{r'=0}^{n-2} \binom{n-2}{r'} p^{r'} q^{(n-2)-r'} \\
 &= n(n-1) p^2 \cdot 1 = n(n-1) p^2 \\
 &= n^2 p^2 - np^2
 \end{aligned}$$

- Such that:

$$\boxed{V(r) = \sum r^2 p_r - n^2 p^2} = np(1-p) = npq$$

Applications to Histograms and Counting Statistics

- Why is this important? Histograms are often Binomially distributed.
- Data either falls A : falls in bin, or \bar{A} : does not fall in bin.
- p = probability of falling in i th bin.
- $+n$ entries, e.g. students in class, histogram = grades.
- Expected number of entries is np .
- Plot of taking distribution several times and checking how many fall in bin i and then plotting that distribution is Binomial.
- Usually you have 1 histogram.
- Look at entries in bin i – n_i/n = fraction of entries in bin i .
- Estimator $p = n_i/n$.
- Expect if you repeated $\Rightarrow n_i$ would follow Binomial distribution with mean $\sigma_i = \sqrt{npq}$ and $V_i = \sigma_i^2 = np(1-p)$.

$$p = \frac{n_i}{n}$$

$$\sigma_i = \sqrt{n_i \left(1 - \frac{n_i}{n}\right)} \approx \sqrt{n_i} \text{ if } n \gg n_i$$

- Notes: we do know the total number n , how often is it in bin i .
- HW: given the distribution, how many times n do I need to do it to get that.
- r fixed n , vs. n fixed r .

7 Thursday, October 2nd 2025

7.1 Review of the Binomial Distribution and Its Properties

- Recall last time:

$$B(r, n, p) = \binom{n}{r} p^r (1-p)^{n-r}$$

$$E(r) = np = \mu$$

$$V(r) = np(1-p) = \sigma^2$$

$$r = \sqrt{np(1-p)}$$

$$p = \frac{\mu}{n}$$

- $\epsilon = \frac{r}{n}$

- Number of detections (people often forget the $(1-p)$ term):

$$n\epsilon \pm \sqrt{n\epsilon(1-\epsilon)}$$

- $\sigma_\epsilon = \frac{1}{n}\sigma_r = \frac{1}{n}\sqrt{r\left(1-\frac{r}{n}\right)} = \frac{1}{\sqrt{n}}\sqrt{\epsilon(1-\epsilon)}$

- Standardized skewness:

$$E\left[\left(\frac{x-\mu}{\sigma}\right)^3\right] = \frac{1-2p}{\sqrt{np(1-p)}}$$

- Excess kurtosis:

$$\frac{1-6p(1-p)}{np(1-p)}$$

7.2 Bernoulli Distribution as a Special Case of the Binomial

- Bernoulli Distribution: Binomial with $n = 1$.

$$B(r, n=1, p) = P_r = \binom{1}{r} p^r (1-p)^{1-r} = p^r (1-p)^{1-r}$$

$$P_0 = 1-p$$

$$P_1 = p$$

$$\mu = E(r) = np = p$$

$$V = \sigma^2 = np(1-p) = p(1-p)$$

$$E(r^k) = \sum_{r=0}^1 r^k P_r = 0^k(1-p) + 1^k p = p$$

- Central moments: $E[(r-p)^k]$

7.3 Negative Binomial and Geometric Distributions

- Negative Binomial: How many n to get r successes.
- Geometric distribution: negative binomial with $r = 1$ (number of trials to get first success).

$$G(n, p) = p(1 - p)^{n-1}$$

- Example: Roll a die with $p(i) = 1/6$ (success = get 4)

$$E(n) = \frac{1}{p} = 6$$

- $P(n \leq 5) = 0.598$
- $P(n \geq 7) = 0.335$

7.4 Samples and the Concept of an Ensemble

- Samples: A set of N draws/trials from a pdf $p(x)$, $\{P_r\}$, is called a **sample** of size N : $\{x_i\}_{i=1}^N$.
- Orthodox statistics: your sample is one of many possible, and we can answer questions about the **ensemble** of samples.

- Samples:

- 1: $\{x_{1i}\}_{i=1}^N$ from pdf $p(x)$
- 2: $\{x_{2i}\}_{i=1}^N$ from pdf $p(x)$
- ...
- M

- $p(x|\mu, \sigma)$
- p = probability of H (heads), $q = 1 - p$ = probability of T (tails)
- These are limits: $p = \lim_{n \rightarrow \infty} \frac{n_H}{n}$
- Flip 10 times: exactly $2^{10} = 1024$ possible outcomes

$$S_1 = \{HHHHHHHHHH\}$$

$$S_2 = \{HTHHHHHHHH\}$$

...

- Given true p , calculate probability of any S_i :

$$P_i = \lim_{N \rightarrow \infty} \frac{n_i}{N}$$

- From the pdf probability, select k samples from a sample space of size F .
- Science/statistics infers from sample space \mathcal{R} to get $p(x)$.
- Not purely deductive: ∞ number of pdfs map to one sample.
- Inductive: If sample pdf is more likely.
- Sample HHTHT ... known, from it we try to infer $p = p_H$.

7.5 Poisson Distribution as a Limit of the Binomial

- Poisson Distribution: limit of binomial for large n , small p .

•

$$B(r, n, p) = \binom{n}{r} p^r (1-p)^{n-r}$$

- Let $n \rightarrow \infty$, $p \rightarrow 0$ such that $np = \mu$ is constant.
- Stirling approximation:

$$n! \approx \sqrt{2\pi n} n^n e^{-n}$$

- Then:

$$B(r, n, p) = \frac{n!}{r!(n-r)!} p^r (1-p)^{n-r}$$

- r is finite, $n \rightarrow \infty$, $n-r \rightarrow \infty$

$$\begin{aligned} B(r, n, p) &= \frac{1}{r!} \frac{\sqrt{2\pi n}}{\sqrt{2\pi(n-r)}} \frac{n^n e^{-n}}{(n-r)^{n-r} e^{-(n-r)}} \left(\frac{\mu}{n}\right)^r \left(1 - \frac{\mu}{n}\right)^{n-r} \\ &= \frac{1}{r!} \sqrt{\frac{n}{n-r}} \left(\frac{n}{n-r}\right)^{n-r} \frac{\mu^r}{e^r} \frac{\left(1 - \frac{\mu}{n}\right)^n}{\left(1 - \frac{\mu}{n}\right)^r} \\ &= \frac{1}{r!} \sqrt{\frac{n}{n-r}} \left(1 - \frac{r}{n}\right)^r \frac{\mu^r}{e^r} \frac{\left(1 - \frac{\mu}{n}\right)^n}{\left(1 - \frac{\mu}{n}\right)^r} \end{aligned}$$

- $n \rightarrow \infty$, $\sqrt{\frac{n}{n-r}} \rightarrow 1$
- $\left(1 - \frac{\mu}{n}\right)^n \rightarrow e^{-\mu}$
- $\left(1 - \frac{\mu}{n}\right)^r \rightarrow 1$
- Then:

$$\lim_{n \rightarrow \infty, np = \mu} B(r, n, p) = \frac{1}{r!} \mu^r e^{-\mu} = P(r|\mu) = \text{Poisson}$$

7.6 Poisson Process and Radioactive Decay

- Consider radioactive decay of some atoms:
 1. Any time interval $[t, t + dt]$ contains at most one decay.
 2. Probability of a decay occurring in this interval is proportional to dt .
 3. Whether or not an atom decays in the interval is independent of any other non-overlapping interval.
- From (1) and (2):

$$P_d(dt) = \lambda dt$$

- Probability of no decay in interval:

$$P_0(dt) = 1 - \lambda dt$$

- Probability of no decay by time $t + dt$:

$$\begin{aligned} P_0(t + dt) &= P_0(t)P_0(dt) \\ &= P_0(t)(1 - \lambda dt) \\ P_0(t + dt) - P_0(t) &= -\lambda P_0(t)dt \\ \frac{dP_0(t)}{dt} &= -\lambda P_0(t) \\ P_0(t) &= P_0(0)e^{-\lambda t} = e^{-\lambda t} \end{aligned}$$

- Probability of getting r decays in time $t + dt$:

$$\begin{aligned} P_r(t + dt) &= P_r(t)P_0(dt) + P_{r-1}(t)P_d(dt) \\ &= P_r(t)(1 - \lambda dt) + P_{r-1}(t)\lambda dt \\ \frac{dP_r(t)}{dt} &= -\lambda P_r(t) + \lambda P_{r-1}(t) \end{aligned}$$

- Solution to PDE:

$$P_r(t) = \frac{1}{r!}(\lambda t)^r e^{-\lambda t}$$

- Poisson distribution with $\mu = \lambda t$

7.7 Moments and Variance of the Poisson Distribution

- Properties:

$$\begin{aligned} E(r) &= \sum_{r=0}^{\infty} r P(r, \mu) \\ &= \sum_{r=0}^{\infty} r \frac{\mu^r}{r!} e^{-\mu} \\ &= \mu e^{-\mu} \sum_{r=1}^{\infty} \frac{\mu^{r-1}}{(r-1)!} \\ &= \mu e^{-\mu} e^{\mu} \\ &= \mu \end{aligned}$$

- $V(r) = E(r^2) - \mu^2$

$$\begin{aligned} E[r(r-1)] &= E(r^2) - \mu \\ &= \sum_{r=2}^{\infty} r(r-1) \frac{\mu^r}{r!} e^{-\mu} \\ &= \mu^2 e^{-\mu} \sum_{r=2}^{\infty} \frac{\mu^{r-2}}{(r-2)!} \\ &= \mu^2 e^{-\mu} e^{\mu} = \mu^2 \end{aligned}$$

- $E(r^2) - \mu = \mu^2$
- $V(r) = E(r^2) - \mu^2 = \mu$
- Binomial: $\mu = np$, $V(r) = np(1-p) = \mu$ for $p \rightarrow 0$ and $n \rightarrow \infty$.

8 Tuesday, October 7th, 2025

8.1 The Gaussian (Normal) Distribution

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

8.2 The Standard Normal Distribution

$$N(0, 1) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

8.3 Example: Circular Symmetry (Darts on a Board)

Distribution of darts:

$$f(x, y) = h(x)k(y)$$

Transform to polar coordinates:

$$g(r, \theta) \approx f(x, y) = h(x)k(y), \quad g(r, \theta) = g(r)$$

$$\frac{\partial g}{\partial \theta} = 0 = \frac{\partial f}{\partial x} \frac{\partial x}{\partial \theta} + \frac{\partial f}{\partial y} \frac{\partial y}{\partial \theta}$$

Coordinates:

$$\begin{aligned} x &= r \cos \theta, & y &= r \sin \theta \\ \frac{\partial x}{\partial \theta} &= -r \sin \theta = -y, & \frac{\partial y}{\partial \theta} &= r \cos \theta = x \end{aligned}$$

Condition:

$$\begin{aligned} 0 &= h'(x)k(y)(-y) + h(x)k'(y)x \\ \frac{h'(x)}{xh(x)} &= \frac{k'(y)}{yk(y)} = a \quad (\text{constant}) \end{aligned}$$

Solutions:

$$\begin{aligned} h(x) &= ce^{ax^2}, & k(y) &= de^{ay^2} \\ f(x, y) &= Ae^{a(x^2+y^2)} = Ae^{ar^2} \approx Ae^{-r^2} \end{aligned}$$

8.4 Expectation Value of a Gaussian

$$E(x) = \mu = \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} xe^{-\frac{(x-\mu)^2}{2\sigma^2}} dx = \mu$$

Useful identity:

$$\int_{-\infty}^{\infty} e^{-ax^2} dx = \sqrt{\frac{\pi}{a}}$$

8.5 Moments of the Gaussian

The n th central moment:

$$\frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} (x - \mu)^n e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx$$

- All odd moments vanish (symmetry about μ).

Define

$$I_0(a) = \int_{-\infty}^{\infty} e^{-ay^2} dy = \sqrt{\frac{\pi}{a}}$$

Differentiation rule:

$$\frac{d^n I_0(a)}{da^n} = (-1)^n \frac{(2n)!}{n!} \frac{I_0(a)}{(2a)^n}$$

8.6 Variance of the Gaussian

Let $y = x - \mu$. Then

$$\begin{aligned} V(y) = V(x) &= \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} y^2 e^{-\frac{y^2}{2\sigma^2}} dy \\ &= -\frac{dI_0(a)}{da} \Big|_{a=\frac{1}{2\sigma^2}} = \sigma^2 \end{aligned}$$

8.7 Kurtosis of the Gaussian

Fourth central moment:

$$\int y^4 e^{-ay^2} dy = \frac{d^2 I_0(a)}{da^2} \frac{1}{4a^2}$$

Evaluates to:

$$E \left[\left(\frac{x - \mu}{\sigma} \right)^4 \right] = 3$$

Thus the Gaussian kurtosis = 3. - Excess kurtosis = 0. - Distributions with > 3 have “fat tails.”

8.8 Poisson Distribution and Gaussian Limit

$$P(r|\lambda) = \frac{1}{r!} \lambda^r e^{-\lambda}$$

For large r , Stirling approximation:

$$r! \approx \sqrt{2\pi r} \left(\frac{r}{e} \right)^r$$

$$\log P(r|\lambda) = -\log(r!) + r \log \lambda - \lambda$$

Expanding around $r \approx \lambda$ leads to Gaussian limit with variance λ :

$$P(r|\lambda) \approx \frac{1}{\sqrt{2\pi\lambda}} e^{-\frac{(r-\lambda)^2}{2\lambda}}$$

8.9 Central Limit Theorem (CLT)

Let x_1, \dots, x_n be independent random variables with mean μ , variance σ^2 . Define sample mean

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

Then as $n \rightarrow \infty$:

$$\bar{x} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$$

8.10 Cumulative Distribution Function of a Gaussian

$$F(x) = \int_{-\infty}^x G(y|\mu, \sigma) dy$$

Define error function:

$$\text{erf}(t) = \frac{2}{\sqrt{\pi}} \int_0^t e^{-y^2} dy, \quad \text{erfc}(t) = 1 - \text{erf}(t) = \frac{2}{\sqrt{\pi}} \int_t^{\infty} e^{-y^2} dy$$

8.11 Gaussian Confidence Intervals

- 1σ : 68.27%
- 2σ : 95.45%
- 3σ : 99.73%
- 5σ : 99.99994%

8.12 Estimators

Given a sample of size n , an *estimator* is any function designed to estimate a property of the true pdf from which the samples were drawn.

9 Thursday, October 9th 2025

9.1 Properties of Estimators

- Consistent: $\lim_{n \rightarrow \infty} \hat{a} = a$
- Unbiased: $\mathbb{E}[a(\hat{x})] = a$
- Efficient: smallest variance of all unbiased estimators

9.2 Example: Measurements and Models

- Let x_1, x_2, \dots, x_n be n measurement points.
- Example applications: number of elements in a histogram bin, position of hits in a detector.
- y_i are the measured values at each x_i , with variances $V(y_i) = \sigma_i^2$.
- Suppose we suspect a model for the histogram shape (e.g. linear background + Gaussian signal):

$$\text{Number of entries} = mx_i + b + Ae^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

- More generally, assume a function $f(x, \vec{\theta})$ with parameters $\vec{\theta}$.

9.3 Least Squares Estimation

- The best estimate for $\vec{\theta}$ is the value that minimizes the chi-squared:

$$\chi^2 = \sum_{i=1}^n \frac{(y_i - f(x_i, \vec{\theta}))^2}{\sigma_i^2}.$$

- Condition for minimization:

$$\frac{\partial \chi^2}{\partial \theta_j} = 0.$$

- Equivalent system of equations:

$$\sum_{i=1}^n \frac{(y_i - f(x_i, \vec{\theta}))}{\sigma_i^2} \frac{\partial f(x_i, \vec{\theta})}{\partial \theta_j} = 0.$$

9.4 Straight Line Fit

- For $f(x_i, \vec{\theta}) = mx_i + b$, with $\vec{\theta} = (m, b)$:

$$\chi^2 = \sum_{i=1}^n \frac{(y_i - (mx_i + b))^2}{\sigma_i^2}.$$

- Normal equations from minimization:

$$\frac{\partial \chi^2}{\partial m} = -2 \sum_{i=1}^n \frac{(y_i - (mx_i + b))}{\sigma_i^2} x_i = 0,$$

$$\frac{\partial \chi^2}{\partial b} = -2 \sum_{i=1}^n \frac{(y_i - (mx_i + b))}{\sigma_i^2} = 0.$$

- Two equations, two unknowns. Can be written in matrix form:

$$\begin{bmatrix} S_{xx} & S_x \\ S_x & S_1 \end{bmatrix} \begin{bmatrix} m \\ b \end{bmatrix} = \begin{bmatrix} S_{xy} \\ S_y \end{bmatrix},$$

where

$$S_{xx} = \sum_i \frac{x_i^2}{\sigma_i^2}, \quad S_x = \sum_i \frac{x_i}{\sigma_i^2}, \quad S_1 = \sum_i \frac{1}{\sigma_i^2},$$

$$S_{xy} = \sum_i \frac{x_i y_i}{\sigma_i^2}, \quad S_y = \sum_i \frac{y_i}{\sigma_i^2}.$$

9.5 Generalized Least Squares with Covariance Matrix

- In general, for non-diagonal covariance matrix V of y_i :

$$\chi^2 = \sum_{i=1}^n \sum_{j=1}^n (y_i - f(x_i, \vec{\theta})) E_{ij} (y_j - f(x_j, \vec{\theta})),$$

where $E = V^{-1}$ is the inverse covariance matrix.

- Linear case: if $f(x_i, \vec{\theta})$ is linear in θ :

$$\vec{f} = A\vec{\theta}.$$

- Then

$$\chi^2 = (\vec{y} - A\vec{\theta})^T V^{-1} (\vec{y} - A\vec{\theta}).$$

- Minimization gives:

$$(A^T V^{-1} A) \vec{\theta} = A^T V^{-1} \vec{y},$$

$$\Rightarrow \vec{\theta} = (A^T V^{-1} A)^{-1} A^T V^{-1} \vec{y}.$$

9.6 Covariance of the Estimated Parameters

- Propagation of covariance:

$$V(\vec{y}) = B V(\vec{x}) B^T.$$

- For parameter estimates:

$$V(\vec{\theta}) = (A^T V^{-1} A)^{-1}.$$

9.7 Goodness of Fit

- The chi-squared statistic

$$\chi^2 = \sum_{i=1}^n \frac{(y_i - f(x_i, \vec{\theta}))^2}{\sigma_i^2}$$

is distributed as χ^2 with n degrees of freedom if the y_i are Gaussian.

- If the model is good, $\chi^2/\text{dof} \sim 1$; if the model is bad, $\chi^2/\text{dof} \gg 1$.

10 Tuesday, October 14th 2025

In this lecture, I start to not take all notes, instead I write down only the key points.

10.1 Chi-Squared for Uncorrelated and Correlated Measurements

- Recall difference in chi-squared formula for uncorrelated and correlated measurements.
- Uncorrelated:

$$\chi^2 = \sum_{i=1}^n \frac{(y_i - \mu_i)^2}{\sigma_i^2}$$

- Correlated (general case):

$$\chi^2 = \sum_{i,j}^n (y_i - \mu_i) V_{ij}^{-1} (y_j - \mu_j)$$

where V^{-1} is the inverse of the covariance matrix V .

10.2 Covariance Matrix and Linear Transformations

- The covariance matrix is defined as

$$V_{ij} = \text{Cov}(y_i, y_j) = \langle (y_i - \mu_i)(y_j - \mu_j) \rangle.$$

- Suppose we apply a linear transformation B to \vec{y} , where B is an $n \times n$ matrix of eigenvectors that diagonalizes V .
- Define transformed variables:

$$\vec{z} = B\vec{y}, \quad V_z = BVB^T, \quad \mu_z = B\mu_y.$$

- Then:

$$\chi^2(z) = (\vec{z} - \vec{\mu}_z)^T V_z^{-1} (\vec{z} - \vec{\mu}_z)$$

10.3 Modeling Data with Parameters

- Suppose we have a model

$$y_i = f(x_i, \vec{\theta}),$$

where x_i are independent variables and $\vec{\theta}$ are model parameters.

- Example: $y_i = mx_i + b$, with $\sigma_i \neq \sigma_j$ (heteroscedastic errors).
- We want to find the best estimate of $\vec{\theta}$.

10.4 Least Squares Estimation

- Define:

$$\chi^2(\vec{\theta}) = \sum_{i=1}^n \frac{(y_i - f(x_i, \vec{\theta}))^2}{\sigma_i^2}.$$

- The best estimators for μ_i are $f(x_i|\vec{\theta})$ and

$$\chi^2(\hat{\theta}) = \chi^2_{\min} = \sum_i \left(\frac{y_i - f(x_i|\hat{\theta})}{\sigma_i} \right)^2.$$

- Principle of least squares:

$$\hat{\theta} = \arg \min \chi^2(\vec{\theta})$$

i.e. the value of $\vec{\theta}$ that minimizes $\chi^2(\vec{\theta})$.

- Solution satisfies:

$$\frac{\partial \chi^2}{\partial \theta_j} = 0.$$

10.5 Distribution of Parameter Estimates

- Note: $\hat{\theta}$ is itself a random variable, with its own probability distribution.
- A different sample $\{x_i, y_i\}$ will lead to a different $\hat{\theta}$.

10.6 Quadratic Expansion of chi-squared and Error Estimates

- For polynomial (linear in parameters) fits:

$$f(x_i|\vec{\theta}) = (A\vec{\theta})_i.$$

- Then:

$$\begin{aligned} \chi^2 &= (\vec{y} - A\vec{\theta})^T V^{-1} (\vec{y} - A\vec{\theta}) \\ \chi^2 &= (y_i - A_{im}\theta_m) V_{ij}^{-1} (y_j - A_{jn}\theta_n). \end{aligned}$$

- First derivative:

$$\frac{\partial \chi^2}{\partial \theta_k} = -2(y_i - A_{im}\theta_m) V_{ij}^{-1} A_{jk}.$$

- Second derivative:

$$\frac{\partial^2 \chi^2}{\partial \theta_k \partial \theta_l} = 2A_{il} V_{ij}^{-1} A_{jk} = 2(A^T V^{-1} A)_{kl}.$$

10.7 Covariance of Parameter Estimates

- For linear fits:

$$V(\hat{\theta}) = (A^T V^{-1} A)^{-1}.$$

- For non-linear fits:

$$V^{-1}(\hat{\theta}) = \frac{1}{2} \frac{\partial^2 \chi^2}{\partial \theta_j \partial \theta_k} \Big|_{\theta=\hat{\theta}}.$$

11 Thursday, October 16th 2025

I was gone this day, so I copied notes from J. Liang.

11.1 Chi-Squared Minimization and Degrees of Freedom

- Degrees of Freedom (dof):

$$\chi^2(\theta) = (y - A\vec{\theta})^T V^{-1} (y - A\vec{\theta})$$

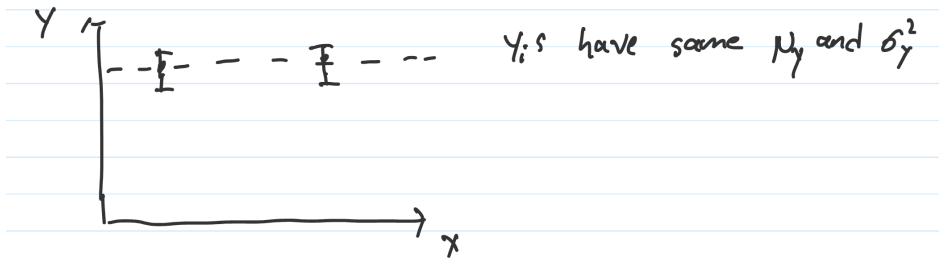
- At $\hat{\theta}$, χ^2 is minimized:

$$\left. \frac{\partial \chi^2}{\partial \theta} \right|_{\theta=\hat{\theta}} = 0 = F(\vec{y}|\vec{\theta}, \vec{x}) = \begin{matrix} F(\vec{y}|\theta_1, \vec{x}) \\ F(\vec{y}|\theta_2, \vec{x}) \\ \vdots \\ F(\vec{y}|\theta_k, \vec{x}) \end{matrix}$$

- i.e. for a linear fit $\vec{y} = A\vec{\theta}$ with k equations:

$$\hat{\theta} = (A^T V^{-1} A)^{-1} (A^T V^{-1}) \vec{y}$$

- Think about it like this: If I know $n - k$ of the y_i 's, the remaining k y_i 's are fixed. Their relations might be complex but they are fixed since we have k equations.



y_i 's have some N_y and σ_y^2 .

11.2 Two-Measurement Example and Correlated Variables

- Best estimator of the true N is $\hat{y} = \frac{y_1 + y_2}{2}$.

$$\chi^2 = \left(\frac{y_1 - \hat{N}}{\sigma_y} \right)^2 + \left(\frac{y_2 - \hat{N}}{\sigma_y} \right)^2$$

- Note:

$$y_1 - \hat{N} = \frac{1}{2}(y_1 - y_2)$$

$$y_2 - \hat{N} = \frac{1}{2}(y_2 - y_1)$$

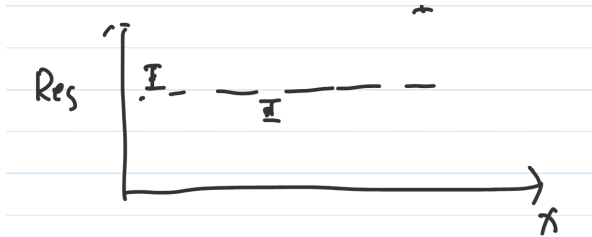
- Where $z_1 = y_1 - y_2$, and $z_2 = y_2 - y_1$ such that $z_2 = -z_1$.

$$V(z_1, z_2) = \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}$$

- Determinant of V is 0, so it is not invertible.
- Thus we effectively have only one independent variable.

11.3 Residuals and Goodness of Fit

- Residual: $r_i = y_i - f_i(\vec{\theta})$.
- When talking about goodness of fit, people usually divide χ^2 by the number of degrees of freedom (dof). If $\chi^2/\text{dof} \approx 1$, it is a good fit.
- What about residuals?



Sometimes one weird data point can throw off the whole χ^2 in unexpected ways, so it is important to check the residuals as well.

- If residuals are randomly scattered around 0, it indicates a good fit.

11.4 Distribution of Estimators

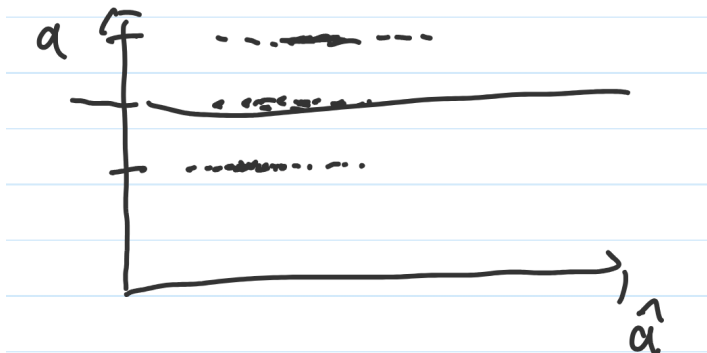
- We have been talking about estimators, but we want values of the parameters.
- Suppose we perform a fit on a parameter with true value a .
- From fitting, we get different estimator values from different data sets.



- These estimates can be scattered over a range of values.

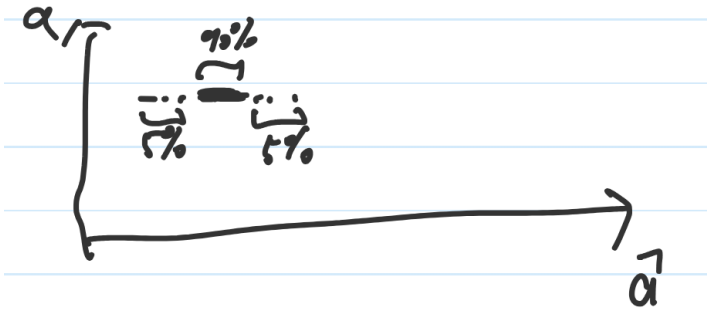
11.5 Toy Monte Carlo Simulations for Estimator Distributions

- Toy Monte Carlo:

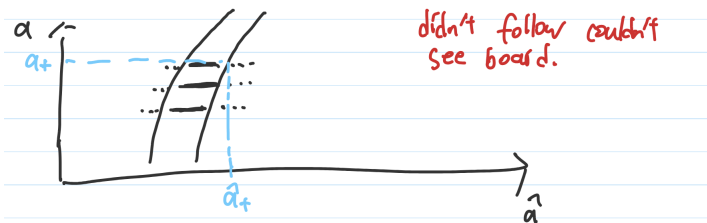


- Often, the model is too complicated to get an analytic form of the estimator distribution.

1. Pick a parameter a .
 2. Generate many data sets according to the model with parameter a .
 3. For each data set, compute the estimator \hat{a} .
 4. Plot a histogram of \hat{a} .
 5. Repeat steps 2–4 n times.
 6. Repeat step 1 m times.
- For each a value, we can determine e.g. 5%, 90%, and 5% quantiles.



- We then connect these points.

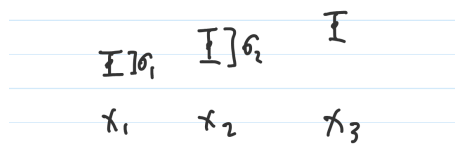


- $P(a_- \leq a \leq a_+) = 0.90$. Coverage = fraction of time your prescription for the estimator interval contains the true value a .
- In other words, we find or choose \hat{a}_- and \hat{a}_+ such that 90% of the time, the true a lies in our interval (corresponding to a_- and a_+).
- The bias in all this is that we can really pick intervals however we want.



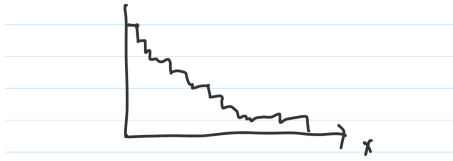
11.6 Typical Applications of Least Squares Fitting

- Typical usage of Least Squares (LS):
 1. x_i, y_i , with known σ_x , and known σ_y for y_i .

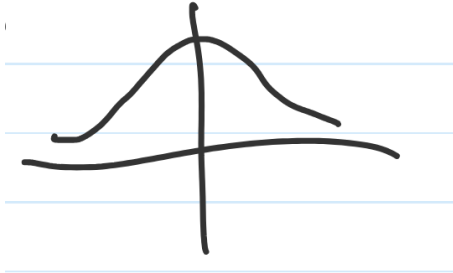


If y_i 's are measured by a detector, we can determine σ_i by analyzing the detector.

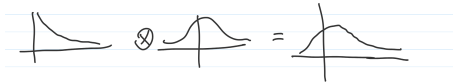
2. Histograms: i.e., measure x on a coordinate, which can have resolution effects, etc., where $y =$ statistics of the number of events at the same x .



Say the resolution has:



so what you see is actually a convolution of the true distribution with the resolution function.



So the fit should actually be an exponential \otimes Gaussian in this case.

$$\chi^2 = \sum_i \frac{(n_i - f_i(\theta))^2}{\sigma_i^2}$$

Suppose we expect Poisson statistics:

$$\sigma_i^2 = n_i$$

Then the data are weighted by n_i :

$$\text{Neyman } \chi^2 \equiv \sum_i \frac{(n_i - f_i(\theta))^2}{n_i}$$

which is a modified least squares form.

Alternatively, it is sensible to say that the expected entries are given by our model:

$$\Rightarrow \text{use } f_i(\vec{\theta}) \text{ as the mean.}$$

Therefore:

$$\text{Pearson } \chi^2 \equiv \sum_i \frac{(n_i - f_i(\theta))^2}{f_i(\theta)}$$

For Neyman χ^2 , if a bin is empty ($n_i = 0$) then it diverges. Both Neyman and Pearson forms are biased in opposite ways. One could use Neyman and Pearson together, but is it worth the effort?

12 Tuesday, October 21st 2025

12.1 Least-Squares Fits

- Usually we have data $(x_i, y_i \pm \sigma_i)$. We want to fit a model $y = f(x; \theta)$ to the data.

$$\chi^2 = \sum_i \frac{(y_i - f(x_i; \theta))^2}{\sigma_i^2}$$

- Sometimes we are given $y = f(x)$ and sometimes we are given $x = g(y)$.

$$\chi^2 = \sum_i \frac{(x_i - g(y_i; \theta))^2}{\sigma_{x,i}^2}$$

$$\sigma_x = \left| \frac{dg}{dy} \right| \sigma_y$$

- Now that is the end of least-squared fits for a while.

12.2 Unbinned Data and Likelihood Functions

- Now we move on to unbinned data.
- Idea is that you have some data drawn from some probability distribution $P(x; a)$

$$P(t) \sim \frac{1}{\tau} e^{-t/\tau}$$

- Data sample of size n : $\{x_1, x_2, \dots, x_n\}$
- Form likelihood function:

$$\mathcal{L}(x_1, x_2, \dots, x_n; a) = \prod_{i=1}^n P(x_i; a)$$

- This is equivalent to the probability of getting the data given the parameter a : $P(\vec{x}|a)$.
- Not a probability distribution in a ! It is a function of a .
- Suppose we have an estimator \hat{a} for a . Then the expectation value of the estimator is:

$$E[\hat{a}] = \int \hat{a}(\vec{x}) P(\vec{x}; a) d\vec{x}$$

- The maximum likelihood principle states that the best estimate for a is the value \hat{a} that maximizes $\mathcal{L}(\vec{x}; a)$:

$$\left. \frac{\partial \mathcal{L}(\vec{x}; a)}{\partial a} \right|_{a=\hat{a}} = 0$$

- Often easier to maximize $\ln \mathcal{L}$ since \ln is monotonic:

$$\ln \mathcal{L}(\vec{x}; a) = \ln \prod_{i=1}^n P(x_i; a) = \sum_{i=1}^n \ln P(x_i; a)$$

- Then we will find:

1. max for $\ln \mathcal{L}(\vec{x}; a)$
 2. min for $-\ln \mathcal{L}(\vec{x}; a)$
- ML estimators tend to be not unbiased, but consistent, often efficient.

$$\left. \frac{\partial \ln \mathcal{L}(\vec{x}; a)}{\partial a} \right|_{a=\hat{a}} = 0$$

12.3 Example: Exponential Distribution

- Example: Exponential distribution

$$P(t; \tau) = \frac{1}{\tau} e^{-t/\tau}$$

$$\mathcal{L}(t_1, t_2, \dots, t_n; \tau) = \prod_{i=1}^n \frac{1}{\tau} e^{-t_i/\tau} = \sum_{i=1}^n \ln \left(\frac{1}{\tau} e^{-t_i/\tau} \right) = -n \ln \tau - \frac{1}{\tau} \sum_{i=1}^n t_i$$

$$\frac{\partial \ln \mathcal{L}}{\partial \tau} = 0 = -\frac{n}{\tau} + \frac{1}{\tau^2} \sum_{i=1}^n t_i$$

$$\hat{\tau} = \frac{1}{n} \sum_{i=1}^n t_i = \bar{t}$$

This is unbiased.

•

$$E(\hat{\tau}) = \frac{1}{n} E\left(\sum_{i=1}^n t_i\right) = \frac{1}{n} \sum_{i=1}^n E(t_i) = \frac{1}{n} n \tau = \tau$$

12.4 Example: Lifetime with Cutoff T

- (Not normalized P):

$$P(t|\tau) = \begin{cases} \frac{1}{\tau} e^{-t/\tau} / e^{-T/\tau} & 0 \leq t \leq T \\ 0 & \text{otherwise} \end{cases}$$

- Now must normalize! $\sum_{t=T}^{\infty} P(t|\tau) = 1$

$$P(t|\tau) = \frac{1}{(1 - e^{-T/\tau})\tau} e^{-t/\tau}$$

- Log-likelihood:

$$\begin{aligned} \ln \mathcal{L} &= \sum_{i=1}^n \left[\ln \left((1 - e^{-T/\tau})\tau \right) \right] - \frac{1}{\tau} \sum_{i=1}^n t_i \\ &= -n \ln(1 - e^{-T/\tau}) - n \ln \tau - \frac{1}{\tau} \sum_{i=1}^n t_i \end{aligned}$$

- Set derivative to zero:

$$\frac{\partial \ln \mathcal{L}}{\partial \tau} = 0$$

$$\frac{\partial \ln \mathcal{L}}{\partial \tau} = \frac{-n(-e^{-T/\tau})(\frac{T}{\tau^2})}{(1 - e^{-T/\tau})} - \frac{n}{\tau} + \frac{1}{\tau^2} \sum_{i=1}^n t_i = 0$$

12.5 Example: Multiple Gaussian Measurements

- Multiple measurements of some quantity:

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

$$\ln \mathcal{L} = -\sum_i \ln(\sqrt{2\pi}) - \sum_i \ln \sigma - \frac{1}{2} \sum_i \frac{(x_i - \mu)^2}{\sigma^2}$$

$$\frac{\partial \ln \mathcal{L}}{\partial \mu} = 0 = \frac{-1}{2} (-2) \sum_i \frac{(x_i - \mu)}{\sigma^2}$$

$$= \frac{1}{\sigma^2} \sum_i (x_i - \hat{\mu}) = 0$$

$$\sum_i x_i = n\hat{\mu} \Rightarrow \hat{\mu} = \frac{1}{n} \sum_i x_i = \bar{x}$$

- Now for σ :

$$\left. \frac{\partial \ln \mathcal{L}}{\partial \sigma} \right|_{\hat{\sigma}, \hat{\mu}} = 0 = -\frac{n}{\sigma} - \frac{1}{2} \left(\frac{-2}{\sigma^3} \right) \sum_i (x_i - \mu)^2$$

$$-n\hat{\sigma}^2 + \sum_i (x_i - \hat{\mu})^2 = 0$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_i (x_i - \hat{\mu})^2 = \frac{1}{n} \sum_i (x_i - \bar{x})^2$$

- This is a biased estimator for σ^2 . Unbiased is with $1/(n-1)$.

12.6 Properties of the Maximum-Likelihood Estimator

- This next stuff is not really testable but can be interesting to see where it comes from.
- To avoid confusion call a_0 the true value of a .
- We have $\mathcal{L}(\vec{x}; a)$ and we want to know how well \hat{a} estimates a_0 .

$$\left. \frac{\partial \ln \mathcal{L}(\vec{x}; a)}{\partial a} \right|_{a=\hat{a}} = 0$$

- Taylor expand around a_0 :

$$f(\hat{a}) = f(a_0) + f'(a_0)(\hat{a} - a_0) + \frac{1}{2} f''(a_0)(\hat{a} - a_0)^2 + \dots$$

- So we have:

$$\left. \frac{\partial \ln \mathcal{L}(\vec{x}; a)}{\partial a} \right|_{a=\hat{a}} = \left. \frac{\partial \ln \mathcal{L}(\vec{x}; a)}{\partial a} \right|_{a=a_0} + (\hat{a} - a_0) \left. \frac{\partial^2 \ln \mathcal{L}(\vec{x}; a)}{\partial a^2} \right|_{a=a_0} + \dots = 0$$

- For $n \rightarrow \infty$, $\hat{a} \rightarrow a_0$ (consistent estimator), so we can neglect higher order terms.

$$\left. \frac{\partial \ln \mathcal{L}(\vec{x}; a)}{\partial a} \right|_{a=a_0} \rightarrow \left. \frac{\partial \ln \mathcal{L}(\vec{x}; a)}{\partial a} \right|_{a=\hat{a}} = 0$$

13 Thursday, October 23rd 2025

13.1 Definition of the Likelihood Function

- Likelihood

$\mathcal{L}(\vec{x}|\theta) = \text{probability of } \vec{x} \text{ and not } \theta$

- Model $\vec{\theta}$ with $p(\vec{x}|\vec{\theta})$

$$\mathcal{L} = \prod_{i=1}^N p(x_i|\vec{\theta})$$

$$\ln \mathcal{L} = \sum_{i=1}^N \ln p(x_i|\vec{\theta})$$

13.2 Maximum Likelihood Estimation (MLE)

- Maximize likelihood to get best estimate of $\vec{\theta}$. Choose $\hat{\theta}$ such that:

$$\hat{\theta} = \operatorname{argmax}_{\vec{\theta}} \mathcal{L}(\vec{x}|\vec{\theta})$$

$$\left. \frac{\partial \ln \mathcal{L}(\vec{\theta})}{\partial \theta} \right|_{\hat{\theta}} = 0$$

13.3 Quadratic Approximation of the Log-Likelihood

- Shape of \vec{a} distribution: $\ln \mathcal{L}(a)$ around \hat{a} is approximately quadratic.
- True $a = \hat{a}$; expand about a_0 :
- Taylor expansion:

$$f(a) = f(a_0) + (\hat{a} - a_0)f'(a_*) \quad \text{where } a_* \text{ is between } \hat{a} \text{ and } a_0$$

- So, for $f = \frac{\partial \ln \mathcal{L}(a)}{\partial a}$:

$$0 = \left. \frac{\partial \ln \mathcal{L}(a)}{\partial a} \right|_{a_0} + (\hat{a} - a_0) \left. \frac{\partial^2 \ln \mathcal{L}(a)}{\partial a^2} \right|_{a_*}$$

13.4 Asymptotic Limit and Expectation Relation

- Large n for consistent $\hat{a} \rightarrow a_0$:

$$\lim_{n \rightarrow \infty} \left. \frac{\partial^2 \ln \mathcal{L}(a)}{\partial a^2} \right|_{a_*} = \lim_{n \rightarrow \infty} \sum_i \left. \frac{\partial^2 \ln p(x_i|a)}{\partial a^2} \right|_{a_*} \approx \lim_n n \int p(x|a) \left. \frac{\partial^2 \ln p(x|a)}{\partial a^2} \right|_{a_*} dx$$

- Sum over samples x_i drawn from $p(x|a)$:

$$\begin{aligned} &= \lim_{n \rightarrow \infty} n E \left(\left. \frac{\partial^2 \ln p(x|a)}{\partial a^2} \right|_{a_*} \right) \\ &= E \left(\left. \frac{\partial^2 \ln \mathcal{L}}{\partial a^2} \right|_{a_*} \right) \end{aligned}$$

$$\hat{a} - a_0 = - \frac{\left. \frac{\partial \ln \mathcal{L}(a)}{\partial a} \right|_{a_0}}{E \left(\left. \frac{\partial^2 \ln \mathcal{L}}{\partial a^2} \right|_{a_*} \right)}$$

$$0 = \left. \frac{\partial \ln \mathcal{L}(a)}{\partial a} \right|_{a_0} + (\hat{a} - a_0) \left. \frac{\partial^2 \ln \mathcal{L}(a)}{\partial a^2} \right|_{a_*}$$

13.5 Normalization of the Likelihood Function

- The likelihood is normalized:

$$\int \mathcal{L}(\vec{x}|a) d\vec{x} = 1$$

$$\Rightarrow \int \frac{\partial \mathcal{L}(\vec{x}|a)}{\partial a} d\vec{x} = 0$$

- Relation between \mathcal{L} and $\ln \mathcal{L}$:

$$\frac{\partial \mathcal{L}}{\partial a} = \frac{\partial \ln \mathcal{L}}{\partial a} \mathcal{L}$$

$$\frac{\partial \ln \mathcal{L}}{\partial a} = \sum_{i=1}^n \frac{\partial \ln p(x_i|a)}{\partial a}$$

13.6 Gaussian Approximation via the Central Limit Theorem

- The sum of n variables with zero mean:
- By the Central Limit Theorem, for large n , $\frac{\partial \ln \mathcal{L}}{\partial a}$ is Gaussian with mean 0.

$$E \left(\frac{\partial^2 \ln \mathcal{L}}{\partial a^2} \right) = -E \left(\left(\frac{\partial \ln \mathcal{L}}{\partial a} \right)^2 \right)$$

13.7 Variance of the Estimator and the Fisher Information

- Variance of $\hat{a} - a_0$:

$$\text{Var}(\hat{a} - a_0) = \frac{\text{Var} \left(\left. \frac{\partial \ln \mathcal{L}(a)}{\partial a} \right|_{a_0} \right)}{\left(E \left(\left. \frac{\partial^2 \ln \mathcal{L}}{\partial a^2} \right|_{a_*} \right) \right)^2}$$

$$= \frac{E \left(\left(\left. \frac{\partial \ln \mathcal{L}(a)}{\partial a} \right|_{a_0} \right)^2 \right)}{\left(E \left(\left. \frac{\partial^2 \ln \mathcal{L}}{\partial a^2} \right|_{a_*} \right) \right)^2}$$

$$= - \frac{E \left(\frac{\partial^2 \ln \mathcal{L}}{\partial a^2} \Big|_{a_0} \right)}{\left(E \left(\frac{\partial^2 \ln \mathcal{L}}{\partial a^2} \Big|_{a_*} \right) \right)^2}$$

- When $n \rightarrow \infty$, $a_* \rightarrow a_0$ and $\hat{a} \rightarrow a_0$:

$$\boxed{\text{Var}(\hat{a}) = - \frac{1}{E \left(\frac{\partial^2 \ln \mathcal{L}}{\partial a^2} \Big|_{a_0} \right)}}$$

- Fisher Information matrix:

$$E \left(\frac{\partial^2 \ln \mathcal{L}}{\partial a^2} \right)$$

- For large n , $\hat{a} \rightarrow a_0$. Estimate $E(\cdot)$ by the observed value:

$$\frac{\partial^2 \ln \mathcal{L}}{\partial a^2} \Big|_{\hat{a}}$$

- So the estimate of variance of $\hat{a} - a_0$ is:

$$\boxed{\text{Var}(\hat{a} - a_0) = - \frac{1}{\frac{\partial^2 \ln \mathcal{L}}{\partial a^2} \Big|_{\hat{a}}}}$$

13.8 Taylor Expansion Near the Maximum Likelihood Estimate

- Taylor expansion again:

$$\frac{\partial \ln \mathcal{L}(a)}{\partial a} = \cancel{\frac{\partial \ln \mathcal{L}(a)}{\partial a} \Big|_{\hat{a}}}^0 + (a - \hat{a}) \frac{\partial^2 \ln \mathcal{L}(a)}{\partial a^2} \Big|_{\hat{a}} + \dots$$

where $\frac{\partial^2 \ln \mathcal{L}(a)}{\partial a^2} \Big|_{\hat{a}} = - \frac{1}{V(\hat{a})}$

$$\frac{-(a - \hat{a})}{V(\hat{a})} + \dots$$

- Note that $V(\hat{a} - a_0) = V(\hat{a})$ because a_0 is constant and does not change the variance (it just shifts the distribution).
- So,

$$\ln \mathcal{L}(a) = \ln \mathcal{L}(\hat{a}) - \frac{1}{2V(\hat{a})} (a - \hat{a})^2$$

$$\mathcal{L} = \mathcal{L}(\hat{a}) \exp \left(- \frac{1}{2} \left(\frac{a - \hat{a}}{\sigma_{\hat{a}}} \right)^2 \right)$$

13.9 Goodness of Fit and the Kolmogorov–Smirnov Test

- Note that the value you get from the maximum likelihood does not give information on how good the fit is—it is just relative to other values of the parameters.
- Kolmogorov–Smirnov test for goodness of fit (KS):
 1. Order data points $\{t_i\}$ such that $t_0 \leq t_1 \leq t_2 \leq \dots \leq t_N$
 2. Form an accumulator F (same model CDF C).
- Metric:

$$\max |F(t_i) - C(t_i)|$$

14 Tuesday, October 28th 2025

14.1 Taylor Expansion of the Log-Likelihood Around the Maximum

- Continuing from last time:

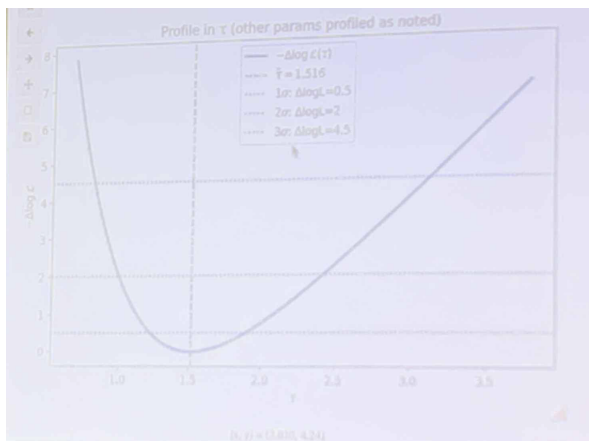
$$\begin{aligned}\ln \mathcal{L}(\vec{\theta}) &= \ln \mathcal{L}(\hat{\theta}) + \sum_i (\theta_i - \hat{\theta}_i) \frac{\partial \ln \mathcal{L}}{\partial \theta_i} \bigg|_{\hat{\theta}} + \frac{1}{2} \sum_{i,j} (\theta_i - \hat{\theta}_i)(\theta_j - \hat{\theta}_j) \frac{\partial^2 \ln \mathcal{L}}{\partial \theta_i \partial \theta_j} \bigg|_{\hat{\theta}} + \dots \\ &= \ln \mathcal{L}(\hat{\theta}) + (\vec{\theta} - \hat{\theta})^T \nabla \ln \mathcal{L} \bigg|_{\hat{\theta}} + \frac{1}{2} (\vec{\theta} - \hat{\theta})^T H (\vec{\theta} - \hat{\theta}) + \dots\end{aligned}$$

- This is the Hessian matrix H :

$$V^{-1} = -H$$

14.2 Parameter of Interest and Nuisance Parameters

- Often interested in only one parameter τ . Pick τ : what is $\max \mathcal{L}$ if that were true.
- τ is what we want. σ is a nuisance parameter.



14.3 Likelihood versus Chi-Squared Interpretation

- Maximum Likelihood is only a function of the shape of your data and whether it matches your model.
- χ^2 is a function of both the height and the shape of your data.

14.4 Normalization and Extended Likelihood

- We can incorporate the number of data events:

$$\int P(x|a) dx = 1 \forall a$$

- Define Q such that:

$$\int Q(x|a)dx = L(a)$$

- Example: number of muons from cosmic rays passing through counters.
- Extend likelihood:

$$\mathcal{L} = \frac{e^{-\nu} \nu^n}{n!}$$

- This is called the extended likelihood and contains the Poisson term.

$$\ln \mathcal{L} = \sum_i \ln(p(x_i|a)) - \nu - n \ln \nu + \ln n! \rightarrow 0$$

14.5 Nuisance Parameters and Systematic Effects

- Nuisance parameters often represent systematic effects or resolution uncertainties:

$$\mathcal{L}(\theta) = \mathcal{L}_0 \frac{1}{\sqrt{2\pi}\sigma_\theta} \exp\left(-\frac{(\theta - \mu)^2}{2\sigma_\theta^2}\right)$$

- $\theta \sim \mathcal{N}(\mu, \sigma_\theta^2)$ is a nuisance parameter.

14.6 Gaussian Data Likelihood and Chi-Squared Minimization

- Suppose we have data points of the form $\{x_i, y_i \pm \sigma_i\}$ and y_i are Gaussian.
- We think we know $y(x) = f(x|a)$.

$$p(y_i|a) = \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left(-\frac{(y_i - f(x_i|a))^2}{2\sigma_i^2}\right)$$

$$\mathcal{L} = \prod_i p(y_i|a)$$

$$\ln \mathcal{L} = -\frac{1}{2} \sum_i \left(\frac{y_i - f(x_i|a)}{\sigma_i}\right)^2 + \text{const}$$

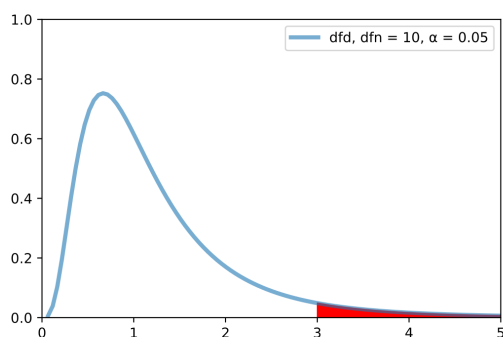
- Maximizing $\ln \mathcal{L}$:
- Minimizing χ^2 :

$$\chi^2 = \sum_i \left(\frac{y_i - f(x_i|a)}{\sigma_i}\right)^2$$

- i.e. Maximizing likelihood is equivalent to minimizing χ^2 for Gaussian data.

14.7 Introduction to Hypothesis Testing

- Next: Hypothesis Testing.
- A simple hypothesis is one for which the test statistic has a completely specified pdf.
- Example: data come from Poisson process with mean $\nu = 5-6$.
- A composite hypothesis is one for which the test statistic has a pdf that depends on unknown parameters.
- Here is a similar example diagram from the internet:



14.8 Rejection Regions and Type I Error

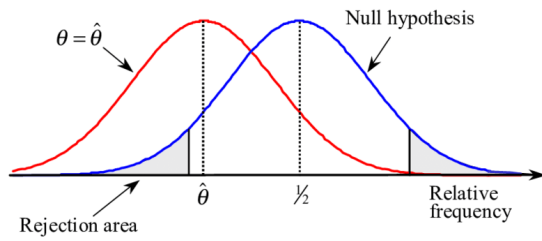
- Decide in advance a σ_i^2 region of x . Reject H_0 if x is in this region.

$$\int_{\text{region}} p(x|H_0)dx = \alpha = \text{significance of test}$$

- For $x \geq x_0$: “Reject” H_0 when true occurs α fraction of the time.
- Type I error: Reject H_0 when true (probability α).
- Fisher.
- $H_0 = \text{Data is } \sim \mathcal{N}(\mu_0, \sigma_0^2)$.
- Get pairs x_1, x_2 .
- Your choice of possible regions is huge and can lead to very weird tests.

14.9 Neyman–Pearson Lemma and Type II Error

- Neyman–Pearson: fixed by adding second hypothesis H_A .



- β = power of test.
- Only two choices: H_0 or H_A .
- Type II error: Accept H_0 when false (probability β).

14.10 Likelihood Ratio Test and Interpretation

- Law: Null hypothesis H_0 is rejected in favor of alternative hypothesis H_A if

$$\Lambda(x) = \frac{\mathcal{L}(x|H_A)}{\mathcal{L}(x|H_0)} \geq k_\alpha$$

- In legal terms: H_0 = innocent until proven guilty, H_A = guilty.
- In science: H_0 = no new physics, H_A = new physics.

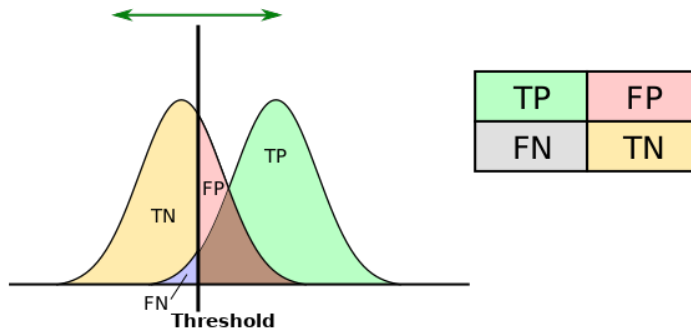
15 Thursday, October 30th 2025

15.1 Introduction to Hypothesis Testing

- Recall the Neyman-Pearson test that we ended last class talking about.
- Hypothesis testing: the test statistic is a function of data.
- Null hypothesis H_0 , alternative hypothesis H_A .
- The main idea from Neyman and Pearson was to include a second hypothesis H_A .

	H_0 true	H_A true
Accept H_0	Correct decision	Type II error
Reject H_0	Type I error	Correct decision

Table 1: Outcomes of hypothesis testing.



Types of errors

		Truth	
		No diff H ₀ to be not rejected	Diff H ₀ to be rejected (H ₁)
Decision based on the p value	H ₀ not rejected No diff	Right decision 1- α	β Type II error
	H ₀ rejected (H ₁) Diff	α Type I error	Right decision 1- β

- H₀ is “true” but rejected: Type I or α error
- H₀ is “false” but not rejected: Type II or β error

17

- Null Hypothesis Significance Testing (NHST).

15.2 The Neyman-Pearson Lemma and Optimal Test Regions

- NP test: How to optimize the region R to minimize Type II error for a fixed Type I error rate α .

$$1 - \beta = \int_R P_A(x) dx$$

- If we move a point x_1 by δx_1 from R to \bar{R} , it must be replaced by another point x_2 such that $\delta_2 P_H(x_2) = \delta_1 P_H(x_1)$. (Gain to loss condition.)

$$\int_R P_H = \alpha$$

- Then change in $1 - \beta$ is $\delta_2 P_A(x_2) - \delta_1 P_A(x_1)$ (gain to loss).
- Select R by picking x values that maximize $P_A(x)/P_{H_0}(x)$ until α is reached:

$$\frac{\delta_2 P_A(x_2)}{\delta_2 P_{H_0}(x_2)} \geq \frac{\delta_1 P_A(x_1)}{\delta_1 P_{H_0}(x_1)}$$

- Decision rule:

$$\frac{P_A(x)}{P_{H_0}(x)} \geq c_\alpha$$

- For some fixed α , set the region by

$$\frac{P_A(x)}{P_{H_0}(x)} = c_\alpha$$

- i.e.

$$\frac{L(\vec{x}|H_A)}{L(\vec{x}|H_0)} \geq c_\alpha$$

15.3 Gaussian Example: Testing Mean Values

- Example: data Gaussian with $\sigma = 1$.

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2}}$$

- Our hypotheses:

$$H_0 : \mu = \mu_0$$

$$H_A : \mu = \mu_1$$

- n measurements:

$$\mathcal{L}(\vec{x}|H_0) = \prod_{i=1}^n f(x_i) = \left(\frac{1}{\sqrt{2\pi}} \right)^n e^{-\sum_{i=1}^n \frac{(x_i - \mu)^2}{2}}$$

- Sample mean:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

- Variance of sample mean:

$$V(\bar{x}) = \frac{\sigma^2}{n} = \frac{1}{n}$$

- Sample variance:

$$S^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

- Expansion:

$$\sum (x_i - \mu_0)^2 = n \left((\bar{x} - \mu_0)^2 + S^2 \right)$$

- Likelihood ratio:

$$\frac{\mathcal{L}(\vec{x}|H_1)}{\mathcal{L}(\vec{x}|H_0)} = e^{-\frac{n}{2}((\bar{x} - \mu_1)^2 - (\bar{x} - \mu_0)^2)} \geq c_\alpha$$

$$\frac{n}{2} \left(2\bar{x}(\mu_1 - \mu_0) + \mu_1^2 - \mu_0^2 \right) \geq \ln(c_\alpha)$$

$$(\mu_1 - \mu_0)\bar{x} \geq \frac{2}{n} \ln(c_\alpha) + \frac{1}{2}(\mu_1^2 - \mu_0^2)$$

- Take the case $\mu_1 > \mu_0$:

$$\bar{x} \geq \frac{1}{\mu_1 - \mu_0} \left(\frac{2}{n} \ln(c_\alpha) + \frac{1}{2}(\mu_1^2 - \mu_0^2) \right)$$

- This is a Simple Statistic (and in our case it is a Sufficient Statistic as well).

- The distribution of \bar{x} is Gaussian about the true mean with width $\sigma/\sqrt{n} = \frac{1}{\sqrt{n}}$.

- So:

$$\int_{\mathcal{R}} \mathcal{L}(\vec{x}|H_0) d\vec{x} = \int_{\bar{x}_\alpha}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{n}{2}(\bar{x} - \mu_0)^2} d\bar{x} = \alpha$$

15.4 Applied Example: Mineral Density Classification

- Example: Mining.
- Opal density = 2.2 g/cm^3 .
- Quartz density = 2.6 g/cm^3 .
- Sample standard deviation $\sigma = 0.2 \text{ g/cm}^3$.
- Gaussian hypotheses:
 - $H_0 \sim N(2.2, \sigma^2)$
 - $H_A \sim N(2.6, \sigma^2)$

- Likelihood ratio:

$$\frac{L(p|H_A)}{L(p|H_0)} = \frac{\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(p-2.6)^2}{2\sigma^2}}}{\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(p-2.2)^2}{2\sigma^2}}} = e^{-\frac{1}{2\sigma^2}((p-2.6)^2 - (p-2.2)^2)} \geq c_\alpha$$

- Cut off at $p \leq 2.56$ (which corresponds to 1.64σ) so we keep 95% of opals.
- But needlessly investigate 36% of useless mines.

15.5 Common Misconceptions in Hypothesis Testing

- A common misconception: if $P(H_0|x)$ is small, then $P(x|H_0)$ must be small. This is **not true**.
- Example: $P(\text{American}|\text{walked on moon}) = 1$, but $P(\text{walked on moon}|\text{American})$ is very small.

16 Tuesday, November 4th 2025

16.1 Likelihood in Hypothesis Testing and Nuisance Parameters

- Likelihood in Hypothesis Testing: Use $\mathcal{L}(\vec{x}|\vec{\theta})$ to select between H_0 (null) and H_1 (alternative). Parameters include both *parameters of interest* and *nuisance parameters*.
- Nuisance Parameters with Constraints: Nuisance parameters often have external measurements that provide probabilistic constraints, e.g., $\theta_{\text{nuisance}} \sim \mathcal{N}(\theta_0, \sigma^2)$. These are incorporated as additional terms in the likelihood:

$$\mathcal{L}_{\text{total}} = \mathcal{L}(\vec{x}|\vec{\theta}) \cdot p(\theta_{\text{nuisance}})$$

16.2 Bayes' Theorem and Its Application to Hypothesis Testing

- Bayes' Theorem: Relates prior knowledge to observations via:

$$P(H|D) = \frac{P(D|H)P(H)}{P(D)}$$

where H = hypothesis, D = data. This allows us to go from $\mathcal{L}(\vec{x}|H)$ to $P(H|\vec{x})$.

- Bayes Theorem approach to hypothesis testing:

$$P(H_0|x) = \frac{P(x|H_0)P(H_0)}{P(x)}$$

- Example: measure c = speed of light. In no sense is c a stochastic variable.
- We have $\mathcal{L}(\vec{x}|H)$ and we want $P(H|\vec{x})$.
- Compare to exponential that when Fourier transformed gives Cauchy distribution.
- Bayes \Rightarrow *posterior* (blue), *likelihood* (red), *prior* (purple), *evidence* (orange).

$$P(H|D) = \frac{P(D|H)P(H)}{P(D)}$$

$$P(H|D) = \frac{\mathcal{L}(D|H)\text{prior}}{\mathcal{L}(D)}$$

$$p(\vec{\theta}|\vec{x}) = \frac{\mathcal{L}(\vec{x}|\vec{\theta})\pi(\vec{\theta})}{\mathcal{L}(\vec{x})}$$

$$\int \mathcal{L}(\vec{x}|\vec{\theta}) d\vec{\theta} = 1$$

- Think conditional probability.

$$p(\vec{x}) = p(\vec{x}|H_0)P(H_0) + p(\vec{x}|H_1)P(H_1)$$

- As long as H_0 and H_1 span all possibilities.

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

- Example: DNA evidence. H_0 = innocent, H_1 = guilty, D = we have DNA evidence at crime.

$$P(D) = P(D|H_0)P(H_0) + P(D|H_1)P(H_1)$$

16.3 Maximum A Posteriori (MAP) Estimation

- Given data of x what value of θ is the most probable: i.e. maximum a posteriori (MAP) estimate (similar to MLE).
- Posterior = $p(\vec{\theta}|\vec{x})$
- Recall: $p(\vec{x}|\vec{\theta}) = \mathcal{L}(\vec{x}|\vec{\theta})$
- Now we have a way to add in prior information about $\vec{\theta}$

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

$$\ln(p(\vec{\theta}|\vec{x})) = \ln(p(\vec{x}|\vec{\theta})) + \ln(p(\vec{\theta})) - \ln(p(\vec{x}))$$

$$\frac{\partial}{\partial \vec{\theta}} \ln(p(\vec{\theta}|\vec{x})) = \frac{\partial}{\partial \vec{\theta}} \ln(p(\vec{x}|\vec{\theta})) + \frac{\partial}{\partial \vec{\theta}} \ln(p(\vec{\theta})) = 0$$

- $\hat{\theta}_{MAP}$: Maximum A Posteriori estimate.
- If prior is flat then MAP = MLE.
- BUT my solution for $\hat{\theta}_{MAP}$ will depend on my choice of prior! i.e. \neq your solution for $\hat{\theta}_{MAP}$
- BUT as data get 'better, more' dependence \sim disappears.
- Example: 99% I am nice, 1% I am evil.
- We can see in demo that with enough data the prior influence goes away, but it takes different amounts of data to finally get to the same answer.

16.4 MAP Estimation for Normal Distribution with Gaussian Prior

- Example: estimate μ of sample. Drawn from $\mathcal{N}(\mu, \sigma^2)$ with σ^2 known.
- Prior, $\mu = \mu_0, \mathcal{N}(\mu_0, \sigma_0^2)$.

$$\mathcal{L} = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right)$$

$$p(\vec{\theta}|\vec{x}) \sim \left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right) \right) \cdot \frac{1}{\sqrt{2\pi\sigma_0}} \exp\left(-\frac{(\mu_0 - \mu)^2}{2\sigma_0^2}\right)$$

$$\log p(\mu|\vec{x}) = -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 - \frac{1}{2\sigma_0^2} (\mu_0 - \mu)^2 + \underbrace{n \log \sqrt{2\pi\sigma}}_{\text{const.}} + \underbrace{\log \sqrt{2\pi\sigma_0}}_{\text{const.}}$$

$$\left. \frac{d \log p(\mu|\vec{x})}{d\mu} \right|_{\hat{\mu}=0} = -\frac{1}{\sigma^2} (-2) \sum_{i=1}^n (x_i - \mu) - 2 \frac{1}{\sigma_0^2} (\hat{\mu} - \mu_0 + 0)$$

$$\frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu) - \frac{\hat{\mu} - \mu_0}{\sigma_0^2} = 0$$

$$\hat{\mu} \left(\frac{n^2}{\sigma^2} + \frac{1}{\sigma_0^2} \right) = \frac{1}{\sigma^2} \sum_{i=1}^n x_i + \frac{1}{\sigma_0^2} \mu_0$$

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

- Rearranging gives:

$$\hat{\mu}_{MAP} = \frac{\frac{1}{\sigma^2} \sum_{i=1}^n x_i + \frac{1}{\sigma_0^2} \mu_0}{\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}}$$

- This is a weighted average of the sample mean and the prior mean.
- As $n \rightarrow \infty$, $\hat{\mu}_{MAP} \rightarrow \hat{\mu}_{MLE}$.

16.5 Conjugate Priors and the Beta Distribution

- Example:

$$P(\vec{\theta}|\vec{x}) = P(\vec{x}|\vec{\theta})P(\vec{\theta})$$

Coin flip: θ = prob. of heads. Bernoulli $\Rightarrow x \in \{0, 1\}$

$$p(x) = \theta^x (1 - \theta)^{1-x}$$

Conjugate priors.

- There used to be an old trick using Beta function.

$$p(\theta) = \text{Beta}(\theta|\alpha, \beta)$$

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

- Now:

$$\log p(\vec{\theta}|\vec{x}) = \log \mathcal{L}(\vec{x}|\vec{\theta}) + \log p(\vec{\theta})$$

- $\hat{\theta}_{MLE}$ is just $\bar{\theta}_{MAP}$ with uniform prior $p(\vec{\theta})$ (i.e. Uninformed prior).

16.6 Choice of Prior: Practical Considerations and Pitfalls

- Example: Proposing Dark Matter (DM) experiment.

- WIMP: $\sim 10^{12}$ eV
- Axion: $\sim 10^{-22}$ eV

- Don't know so take flat prior, but then we end up with massive probability for WIMP and tiny for axion because axion is unbiased.

17 Thursday, November 6th 2025

- The following notes from this entire lecture are not testable.

17.1 Deductive and Inductive Reasoning; Weak Syllogism and Bayesian Motivation

- How do we reason in this realm of deductive vs inductive logic (following textbook from ET Jaynes, from math of Polya around 1450)?
- Aristotelian logic: deductive reasoning from true premises to true conclusions.

$$A \Rightarrow B$$

i.e. if A is true, then B is true.

- You can also say that B is false if A is false.
- But if B is true, A may be true or false.
- Simple Math Example: $A : x > 7, B : x > 5$:
- If $x > 7$ then $x > 5$. i.e. if A true then B true.
- But if A is false (i.e. $x \leq 7$) then B may be true or false.
- or another example: A : I tried to do a “hard-flip” on my skateboard. B : I will get hurt. If not A then B is false.
- Inductive reasoning: from observations to general rules.
- Almost all of science (IRL) reasoning is inductive.
- “weak syllogism”: outcomes are not a certainty, only a likelihood: If B is true, then A becomes more plausible.
- This is the basis of Bayesian reasoning.
- Example: A : If it is raining, B : the ground is wet. Then if the ground is wet. Therefore it is more plausible that it is raining.
- Example: It will start to rain at 10am. B : it will get cloudy before 10am. If we see clouds at 9am, then it is more plausible that it will rain at 10am.
- Example: Police officer at night. A : Hears a burglar alarm. B : sees a jewelry store with broken windows. Sees person in mask climbing out window with a bag of jewelry.
- Science Example: A : If GR is true. B : then perihelion advance of Mercury. If we observe perihelion advance of Mercury, then GR becomes more plausible.
- Science Example: If simplistic SUSY breaking mechanism is true, B : particular type of dark matter.

17.2 Plausibility Notation and Consistency Requirements

- $A|B \equiv$ plausibility of A given B.
- $A|B \in \mathcal{R}$ we will insist on self-consistency of our operations.
- If $A|C$ is “more plausible” than $B|C$, then we write:

$$A|C > B|C$$

- If conclusions can be reached multiple ways, these must give same plausibility.
- Equivalent states of knowledge must be represented by same plausibility.
- Also require Qualitative correspondence with common sense reasoning.
- $C + C'$ contains new information from C such that A becomes more plausible i.e.

$$A|C' > A|C$$

and plausibility of B given A is not changed:

$$B|AC' = B|AC$$

then plausibility of $AB|C' \geq AB|C$.

- Also if $A|C' > A|C$ then $\bar{A}|C' < \bar{A}|C$.
- We then can do some heavy math and show that the only consistent way to assign plausibilities is to identify them with probabilities that follow the rules of probability theory.

17.3 Product Rule for Combined Plausibility

- We try to find a consistent rule for calculating plausibility of A B. We want $AB|C$ (i.e. a product rule).

- Decide $B = \text{True}$, $B|C$
- Given B true decide if A is true: $A|BC$

- Decide $A = \text{True}$, $A|C$, given $A = \text{True}$ we need to decide if $B = \text{T}$, $B|AC$
- Want: $AB|C = ?$
- With $(AB|C)$ we might have functions: $F(B|C, A|BC) = F(A|C, B|AC)$.
- $A|C, B|C, A|BC, B|AC$
- By looking at edge cases:

$$B = \bar{A}$$

- If $C \rightarrow C'$, then $B|C' > B|C$. But plausibility for A given is unchanged:

$$A|BC' = A|BC$$

- Now calculate consistency for: $ABC|D$

- Take triple product and expand into singles as products.

$$\begin{aligned}
 ABC|D &= F(BC|D, A|BCD) \\
 &= F(F(C|D, B|CD), A|BCD) \\
 &= F(C|D, AB|CD) \\
 &= F(C|D)F(B|CD, A|BCD)
 \end{aligned}$$

- Note that for simplicity we can start using new parameters
- From this we get the functional equation:

$$F(F(x, y), z) = F(x, F(y, z))$$

- Associativity equation.

17.4 Solution of the Associativity Equation via Transform

- Can find most general solution:

$$u \equiv F(x, y), \quad v \equiv F(y, z)$$

with x,y,z independent.

$$\boxed{F(x, v) = F(u, z)}$$

- Take derivatives wrt x,y,z.

$$\frac{\partial F(x, v)}{\partial x} \equiv F(x, v) = \frac{\partial F}{\partial u} \frac{\partial u}{\partial x}$$

- $F_1(x, y) = \frac{\partial F(x, y)}{\partial x}$, $F_2(x, y) = \frac{\partial F(x, y)}{\partial y}$

- Divide.

$$\frac{F_2(x, v)}{F_1(x, v)} F_1(y, z) = \frac{F_2(x, y)}{F_1(x, y)}$$

- Set $G(x, y) = \frac{F_2(x, y)}{F_1(x, y)}$

- Then:

1. u:

$$G(x, v)F_1(y, z) = G(x, y)$$

2. v:

$$G(x, v)F_2(y, z) = G(x, y)G(y, z)$$

This last term is not a function of y.

- So:

$$\frac{\partial u}{\partial z} = \frac{\partial v}{\partial y}$$

$$G(x, v)F_{12}(y, z) = 0$$

$$\frac{\partial v}{\partial y} = G(x, v)F_{21}(y, z)$$

$$G(x, y)G(y, z) \text{ independent of } y$$

- Most generic form of $G(x, y) = H(x)/H(y)$

- $w(x) = \exp\left(\int^x H(x')dx'\right)$

- $w(F(y, z)) = w(v) = w(y)w(z)$

18 Thursday, November 13th 2025

- Continuation of the previous lecture on building up the idea of plausibility and converting it into a numerical measure.
- Again this is not testable material.

18.1 Plausibility as a Generalization of Logic

- Plausibility is a generalization of logic to situations where we do not have certainty.

$$(AB|C) \in \mathcal{R}$$

$$(AB|C) = F[(B|C), (A|BC)]$$

$$= F[(A|C), (B|AC)]$$

18.2 Derivation of the Functional Form via Associativity

- Then $(ABC|D)$:

$$F[F(C|D, B|CD), A|BCD] = F[C|D, F(B|CD, A|BCD)]$$

$$F[F(x, y), z] = F[x, F(y, z)]$$

$$u = F(x, y), \quad v = F(y, z)$$

$F(u, z) = F(x, v)$

- Take derivatives wrt x, y, z .

$$F_1(x, y) = \frac{\partial F(x, y)}{\partial x}, \quad F_2(x, y) = \frac{\partial F(x, y)}{\partial y}$$

$$G(x, y) = \frac{F_2(x, y)}{F_1(x, y)}$$

- Then:

$$Q = G(x, v)F_1(y, z) = G(x, y)$$

$$R = G(x, v)F_2(y, z) = G(x, y)G(y, z)$$

- Now use Q and R :

$$\frac{\partial Q}{\partial z} = G_2(x, v)F_2(y, z)F_1(y, z) + G(x, v)F_{12}(y, z) = 0$$

$$\frac{\partial R}{\partial y} = G_2(x, v)F_1(y, z)F_2(y, z) + G(x, v)F_{21}(y, z) = 0$$

- This is used to show that $G(x, y)G(y, z)$ is independent of y .

$$\frac{\partial Q}{\partial y} = 0 \Rightarrow G(x, v)F_{12}(y, z) = 0$$

$$\frac{\partial R}{\partial y} = 0 \Rightarrow G(x, y)G(y, z) \text{ is independent of } y$$

- General form is $G(x, y) = r \cdot H(x)/H(y) \quad r \in \mathcal{R}$

- Pick y_0 s.t. $G(x, y_0) \neq 0, r = G(y_0, y_0)$:
- Define: $H(x) = G(x, y_0)$

$$G(x, y)G(y, z) = G(z, y_0)G(y_0, z)$$

$z = y_0$ gives:

$$G(x, y)G(y, y_0) = G(x, y_0)G(y_0, y_0)$$

$$G(x, y)H(y) = H(x)r$$

$$\boxed{G(x, y) = r \frac{H(x)}{H(y)}}$$

- using our first expression

$$G(x, v)F_1(y, z) = G(x, y)$$

- we have

$$r \frac{H(x)}{H(v)} F_1(y, z) = r \frac{H(x)}{H(y)}$$

$$F_1(y, z) = \frac{H(v)}{H(y)}$$

- then using our second expression

$$r \frac{H(x)}{H(v)} F_2(y, z) = r \frac{H(x)}{H(y)} \cdot r \frac{H(y)}{H(z)}$$

- we have

$$F_2(y, z) = \frac{rH(v)}{H(z)}$$

- Recall that $v = F(y, z)$

$$dv = dF(y, z) = F_1(y, z)dy + F_2(y, z)dz$$

$$dv = \frac{H(v)}{H(y)}dy + \frac{rH(v)}{H(z)}dz$$

- Divide by $H(v)$

$$\frac{1}{H(v)}dv = \frac{1}{H(y)}dy + \frac{r}{H(z)}dz$$

- F was shown earlier to be monotonic, so the derivatives of the functions are non-zero. So G is positive, and H is G's at some fixed point, so H is also non-zero.
- We can then integrate on both sides.

$$\int \frac{1}{H(v)}dv = \int \frac{1}{H(y)}dy + \int \frac{r}{H(z)}dz$$

- Let $I(x) = \int \frac{1}{H(x')}dx'$

$$I(v) = I(y) + I(rz)$$

- We can exponentiate both sides.

$$e^{I(v)} = e^{I(y)} \cdot e^{I(rz)}$$

- Let $w = e^I = \exp \int^x \frac{1}{H(x')} dx'$

$$w(F(y, z)) = w(y) \cdot w^r(z)$$

- Now recall, $u = F(x, y), \quad v = F(y, z)$:

$$F(x, v) = F(u, z)$$

- So:

$$w(F(x, v)) = w(x)w^r(v) = w(F(u, z)) = w^r(z)w(u)$$

- So:

$$r^2 = r \Rightarrow r = 1$$

- So we have:

$$w(F(x, y)) = w(x) \cdot w(y)$$

$$F(x, y) = w^{-1}(w(x) \cdot w(y))$$

18.3 Product Rule for Plausibility

- Recall we wanted the product rule for plausibility.

$$w(AB|C) = w(A|C) \cdot w(B|C)$$

$$w(AB|C) = w(B|AC) \cdot w(A|C)$$

- Suppose A is certain if C is certain (i.e. $C = \text{True}$)

- Then: $(AB|C) = (B|C)$

- Then: $w(AB|C) = w(B|C)$

- Then: $w(AB|C) = w(A|BC) \cdot w(B|C)$

- Then: $w(AB|C) = w(A|C) \cdot w(B|C)$

- So:

$$w(A|C) = 1$$

if A is certain if C is given as True

- If A is impossible given C is true:

$$(A|C) = \text{False}$$

$$(AB|C) = \text{False}$$

- Then:

$$w(AB|C) = w(A|BC) \cdot w(B|C)$$

$$w(A|C) = w(A|C) \cdot w(B|C)$$

- So case $w(A|C) = 0$ (or $w = \infty$), take $w' = 1/w$.

- $\exists w$ satisfies product rule. $0 \leq w(x) \leq 1, \quad \forall$ plausibilities.

18.4 Sum Rule for Plausibility

- Sum rule, any proposition A:

$$A\bar{A} = F$$

$$A + \bar{A} = \text{True}$$

- Set $u = w(A|B)$, $v = w(\bar{A}|B)$
- Must exist S, s.t.

$$v = S(u)$$

- with constraints:

$$S(0) = 1, \quad S(1) = 0$$

- Can swap A and \bar{A} :

$$u \text{ must } = S(v)$$

- So:

$$u = S(S(u))$$

$$S^{-1}(u) = S(u)$$

- Skipping some steps, after lots of algebra we get:

$$S(x) = (1 - x^m)^{1/m}$$

- for $0 \leq x \leq 1$
- for $0 < m < \infty$
- So we have:

$$w(\bar{A}|B) = (1 - (w(A|B))^m)^{1/m}$$

- This is the sum rule for plausibilities.

$$w^m(\bar{A}|B) = 1 - w^m(A|B)$$

$$\boxed{w^m(A|B) + w^m(\bar{A}|B) = 1}$$

For some m.

$$w^m(AB|C) = w^m(A|BC) \cdot w^m(B|C)$$

- Then $w' = w^m$ so take $m = 1$ for simplicity.

18.5 Identification of Plausibility with Probability

- Finally \exists function $w^m(x) \equiv p(x)$ s.t. :

$$p(AB|C) = p(A|BC) \cdot p(B|C)$$

$$= p(B|AC) \cdot p(A|C)$$

$$p(A|C) + p(\bar{A}|C) = 1$$

where $0 \leq p(x) \leq 1$ in monotonic in plausibility.

- Any logical function can be expressed in terms of AND, NOT.

18.6 Application to Deductive Logic

- Deductive logic:

$$A \Rightarrow B$$

- A is True, so B is True.
- B is False, so A is False.
- Define $C = \text{prop that } A \Rightarrow B$

$$p(AB|C) = p(B|AC) \cdot p(A|C)$$

or

$$p(B|AC) = \frac{p(AB|C)}{p(A|C)}$$

and

$$p(A\bar{B}|C) = p(A|\bar{B}C) \cdot p(\bar{B}|C)$$

\Rightarrow

$$p(A|\bar{B}C) = \frac{p(A\bar{B}|C)}{p(\bar{B}|C)}$$

- But $A \Rightarrow B$ means $A\bar{B}$ is impossible.

$$p(A\bar{B}|C) = 0$$

$$p(A|\bar{B}C) = 0$$

- Now another case:

$$\begin{aligned} p(AB|C) &= p(A|BC) \cdot p(B|C) \\ &= p(B|AC) \cdot p(A|C) \end{aligned}$$

\Rightarrow

$$p(A|BC) = \frac{p(B|AC) \cdot p(A|C)}{p(B|C)}$$

$A \Rightarrow B$ means if B is True, then A is more plausible.

$$p(B|AC) = 1$$

$$p(B|C) \leq 1$$

$$\boxed{p(A|C) \leq \frac{p(A|C)}{p(B|C)} = p(A|BC)}$$

- Now $A \Rightarrow B$ means if A is False, then B is less plausible.

$$p(B|\bar{A}C) = \frac{p(\bar{A}|\bar{B}C)p(B|C)}{p(\bar{A}|C)}$$

$$p(A|BC) \geq p(A|C)$$

so $p(\bar{A}|\bar{B}C) \leq p(\bar{A}|C)$

$$\boxed{p(B|\bar{A}C) \leq \frac{p(\bar{A}|C)p(B|C)}{p(\bar{A}|C)} \leq p(B|C)}$$

18.7 Mutually Exclusive and Exhaustive Propositions

- So we will see next class:
- A_1, A_2, \dots, A_n mutually exclusive, exhaustive propositions.
- We also have B proposition.
- A_i = everything that can happen.
- $\Rightarrow p(A_i|B)$ B says nothing
- $A_i A_j = F$
- $p(A_i|B) = \frac{1}{n}$

$$p\left(\sum(A_i)|B\right) = 1$$

- (i.e. $p(\text{certain}|B) = 1$)

19 Tuesday, November 18th 2025

19.1 Continuous Variables: CDF and PDF

- (1)

$$0 \leq p(\cdot) \leq 1$$

- Continuous variables in $y = mx + b$
- f is continuous and real-valued
- $F = f \leq q$
- $F' = f > q$
- given data D

$$p(F|DI) = G(q)$$

- where I is other information
- Take range

- $A \equiv f \leq a$
- $B \equiv f \leq b$
- $C \equiv a \leq f \leq c$

$$B = A + C$$

$$P(AC|I) = 0$$

$$P(B|DI) = P(A|DI) + P(C|DI)$$

- $G(b) = G(a) + P(C|DI)$
- $P(a < f \leq b|DI) = G(b) - G(a)$
- $P(a < f \leq b|DI) = \int_a^b g(f) df$
- where $g(f)$ is the probability density function of f
- $g(f) = \frac{dG}{df}$

19.2 Priors and Conjugate Families

- Priors:
- Conjugate Priors is matched to particular probability distribution forms.
- Example: Beta prior + Bernoulli, Binomial, Geometric,
- Data + Prior \rightarrow Posterior = same family as prior = Beta
- Pressure to pick prior that is least informative, most unbiased.

19.3 Fisher Information and Reparameterization

- Recall from deriving the variance and MLE:

$$\frac{d\mathcal{L}(x|\theta)}{d\theta} = \frac{d \log \mathcal{L}}{d\theta} \mathcal{L}(\vec{x}|\theta)$$

$$1 = \int \mathcal{L}(\vec{x}|\theta) d\vec{x} \quad \forall \theta$$

$$0 = \int \frac{d \log \mathcal{L}}{d\theta} \mathcal{L}(\vec{x}|\theta) d\vec{x} = E_x \left(\frac{d \log \mathcal{L}}{d\theta} \right)$$

- Define Fisher Information

$$0 = \int \left(\frac{d^2 \log \mathcal{L}}{d\theta^2} \mathcal{L} + \left(\frac{d \log \mathcal{L}}{d\theta} \right)^2 \mathcal{L} \right) d\vec{x}$$

$$-E \left(\frac{d^2 \log \mathcal{L}}{d\theta^2} \right) = E \left(\left(\frac{d \log \mathcal{L}}{d\theta} \right)^2 \right) \equiv I(\theta) = \text{Fisher Information}$$

- Multiple vars θ

$$I_{i,j}(\theta) = E \left(\frac{d^2 \log \mathcal{L}}{d\theta_i d\theta_j} \right)$$

- MLE $\hat{\theta}$ Hessian

- + (Co)Variance $V_{i,j}(\hat{\theta}) = -\frac{1}{H_{i,j}(\hat{\theta})} = \frac{1}{I_{i,j}(\hat{\theta})}$

- I is large \Rightarrow large curvature

- Likelihood \mathcal{L} is reparametrization invariant, Hessian (so I) is not.

- Doing this in 1D for simplicity

- $\alpha = g(\theta)$

- $\theta = g^{-1}(\alpha) = h(\alpha)$

- $l(\theta) = \log \mathcal{L}(x|\theta)$

- $\hat{l}(\alpha) = l(h(\alpha))$

$$\frac{d\hat{l}}{d\alpha} = l'(h(\alpha))h'(\alpha)$$

- At MLE point $\hat{\theta}$,

- $a = g(\hat{\theta})$

- $\hat{l}'(g(\hat{\theta})) = l'(\hat{\theta})h'(g(\hat{\theta})) = 0$

- $\hat{\alpha} = \text{MLE in terms of } \alpha = g(\hat{\theta})$

- At MLE:

$$\hat{l}''(\alpha) = l''(h(\alpha))(h'(\alpha))^2 + l'(h(\alpha))h''(\alpha) = l''(h(\alpha))(h'(\alpha))^2$$

$$\hat{l}''(\alpha) = l''(\hat{\theta})h'(\hat{\alpha})^2$$

- Where the $l''(\hat{\theta})$ is the Hessian / curvature in terms of θ
- So:

$$\boxed{\frac{d^2 \log \mathcal{L}}{d\alpha^2} = \frac{d^2 \log \mathcal{L}}{d\theta^2} \left(\frac{d\theta}{d\alpha} \right)^2}$$

- $V(\hat{\alpha}) = V(\hat{\theta}) \left(\frac{d\alpha}{d\theta} \right) \Big|_{\hat{\theta}}$
- If you look at delta log likelihood from MLE point, it is invariant under reparametrization.
- Now:

$$E_x \left(\frac{d \log \mathcal{L}(\vec{x}|\alpha)}{d\alpha} \right) = E_x \left(\frac{d \log \mathcal{L}(\vec{x}|\theta)}{d\theta} \left(\frac{d\theta}{d\alpha} \right)^2 \right)$$

- So Fisher Information transforms as:

$$\boxed{I_{\alpha}(\alpha) = I_{\theta}(\theta) \left(\frac{d\theta}{d\alpha} \right)^2}$$

- Let's let $\pi(\theta)$ be the prior in terms of θ

$$\pi = \text{prior} = \pi_{\theta}(\theta)$$

- π_{α} is another prior in terms of α
- We want:

$$\pi_{\theta}(\theta) d\theta = \pi_{\alpha}(\alpha) d\alpha$$

- We want:

$$\begin{aligned} \pi_{\alpha}(\alpha) &= \pi_{\theta}(\theta) \left| \frac{d\theta}{d\alpha} \right| \\ I_{\alpha}(\alpha) &= I_{\theta}(\theta) \left(\frac{d\theta}{d\alpha} \right)^2 \\ \sqrt{I_{\alpha}(\alpha)} &= \sqrt{I_{\theta}(\theta)} \left| \frac{d\theta}{d\alpha} \right| \end{aligned}$$

- So if we pick:

$$\boxed{\pi = \sqrt{I}}$$

- then this gives us a parameter invariant prior called the Jeffreys Prior.

19.4 Jeffreys Prior: Bernoulli Trial Example

- Example with Bernoulli 1 trial:

$$\mathcal{L}(x|p) = p^x(1-p)^{1-x} \quad x \in \{0,1\}$$

$$\log \mathcal{L} = x \log p + (1-x) \log(1-p)$$

$$\frac{d \log \mathcal{L}}{dp} = \frac{x}{p} - \frac{1-x}{1-p} = \frac{x - x/p - p + x/p}{p(1-p)} = \frac{x}{p(1-p)} - \frac{1}{1-p} = \frac{x}{p(1-p)} - \frac{1}{1-p}$$

$$\frac{d^2 \log \mathcal{L}}{dp^2} = -\frac{x}{p^2(1-p)^2} - \frac{1}{(1-p)^2}$$

- Moving to expectation value:

$$E(x) = p \quad \text{For Bernoulli 1 trial}$$

$$I(p) = -E\left(\frac{d^2 \log \mathcal{L}}{dp^2}\right) = \frac{p(1-2p)}{p^2(1-p)^2} - \frac{1}{(1-p)^2} = \frac{1-2p+p}{p(1-p)^2} = \frac{1}{p(1-p)}$$

- Jeffreys Prior for Bernoulli trial:

$$\pi(p) = \sqrt{I(p)} = \frac{1}{\sqrt{p(1-p)}} = \frac{1}{p^{1/2}(1-p)^{1/2}}$$

$$\pi_p = \text{Beta}(1/2, 1/2)$$

$$\text{Beta}(a, b) = \frac{p^{a-1}(1-p)^{b-1}}{B(a, b)}$$

19.5 Maximum Entropy Principle

- Another way to get priors is the Principle of Maximum Entropy (i.e. least information)
- Shannon showed can quantify information using entropy.

$$p_i = \text{prob of outcome } x_i$$

$$H(p) = -\sum_i p_i \log p_i = \text{information}$$

- Rule: least information consistent with whatever else I know about the data.
- Use Lagrange multipliers

$$L(\vec{p}, \lambda) = -\sum_i p_i \log p_i + \lambda \left(\sum_i p_i - 1 \right)$$

$$\frac{\partial L}{\partial p_i} = -\log p_i - 1 + \lambda$$

$$\log p_i = \lambda - 1$$

$$p_i = p_j \quad \forall i, j$$

$$\sum_i p_i = 1 \Rightarrow p_i = \frac{1}{n}$$

- So maximum entropy prior with no other constraints is uniform.

$$np_i = 1$$

- we might know $\mu = \text{mean}$, so then you would have:

$$\lambda_1(\sum_i p_i - 1) + \lambda_2(\sum_i x_i p_i - \mu)$$

- Then you would get an exponential distribution. We will show this more in the next lecture.
- Another example is if we know variance, then we get a Gaussian distribution.

20 Thursday, November 20th 2025

20.1 Overview of Prior Distributions

- Priors:
- Jeffreys Prior:

$$\pi(\theta) = \sqrt{I(\theta)}$$

- Reference Prior:

$$\pi(\vec{\theta}) = \sqrt{\det I(\vec{\theta})}$$

- Max Entropy Prior – knowledge of future e.g. mean
- Others...

20.2 Maximum Entropy Prior with Known Mean

- Constraint:

$$\sum_i p_i = 1$$

- Max entropy:

$$p_i = 1/n$$

- know mean $\mu = \sum x_i p_i$
- using Lagrange multipliers form:

$$L(p, \lambda_0, \lambda_1) = - \sum_i p_i \log p_i + \lambda_0 \left(\sum_i p_i - 1 \right) + \lambda_1 \left(\sum_i x_i p_i - \mu \right)$$

- Take all derivatives and solve.

$$\frac{\partial L}{\partial p_i} = -(\log p_i + 1) + \lambda_0 + \lambda_1 x_i = 0$$

$$\log p_j = \lambda_0 - 1 + \lambda_1 x_j$$

$$p_j = \exp(\lambda_0 - 1) \exp(\lambda_1 x_j)$$

$$p_j = C \exp(\beta x_j)$$

$$\sum_i p_i = 1 \Rightarrow C \sum e^{\beta x_i} = 1$$

$$C = \frac{1}{\sum_i e^{\beta x_i}} = \frac{1}{Z(\beta)}$$

- Know the mean means that your prior has an exponential form.

- Let us see:

$$p_i = e^{\beta x_i} / Z(\beta)$$

with μ constraint

$$\mu = \frac{\sum x_i e^{\beta x_i}}{Z(\beta)}$$

$$\frac{\partial Z}{\partial \beta} = \sum_i x_i e^{\beta x_i} = \frac{Z'(\beta)}{Z(\beta)} \equiv \mu(\beta) = \frac{\partial \log Z}{\partial \beta}$$

- Form $\mu(\beta) - \mu$, Find root β such that $\mu(\beta) - \mu = 0$
- Use numerical methods.

20.3 Maximum Entropy for Continuous Variables

- Often seen as generalization of entropy definition to continuous variables:

$$\sum p_i \log p_i \rightarrow H(\theta) = - \int p(\theta) \log p(\theta) d\theta$$

- Start with binned distribution:

$$H_{\Delta} = - \sum_i p_i \log p_i$$

with width Δx_i

$$p_i = \int_{\text{ith bin}} p(x) dx$$

$$H_{\Delta} = - \sum p(x_i) \Delta x_i \log(p(x_i) \Delta x_i)$$

- Take same width $\Delta x_i = \Delta$

$$= - \sum p(x_i) \Delta \log p(x_i) - \log(\Delta) \left(\sum p(x_i) \Delta \right)$$

$$\rightarrow = - \int p(x) \log(p(x)) dx - \log(\Delta) \int p(x) dx$$

where the log delta term goes to infinity as $\Delta \rightarrow 0$ and integral of $p(x)$ is 1.

- Max entropy of continuous variable in range $[a, b]$:

$$H(\theta) = - \int_a^b p(\theta) \log p(\theta) d\theta + \lambda \left(\int_a^b p(\theta) d\theta - 1 \right)$$

- Then:

$$\frac{\delta H}{\delta p(\theta)} = - (\log p(\theta) + 1) + \lambda = 0$$

$$\log p(\theta) = \lambda - 1$$

$$p(\theta) = \text{Constant} = C$$

- Apply constraint:

$$\int_a^b C d\theta = 1$$

$$C = \frac{1}{b-a}$$

i.e. $p = 1/(b-a)$ uniform distribution.

- If $a, b = (-\infty, \infty)$ then no max entropy distribution exists.
- This is called the “Improper Prior”, $\rightarrow \int \neq 1$
- Then:

$$p(\theta|D) = \frac{p(D|\theta)p(\theta)}{\int p(D|\theta)p(\theta) d\theta}$$

20.4 Scale Invariant and Location Invariant Priors

- Measure data in units
- get x
- Switch to new units in scale factor α

$$x' = \alpha x$$

$$p(x|I) dx = p(x'|I) dx'$$

$$dx' = \alpha dx$$

$$p(x|I) dx = p(\alpha x|I) \alpha dx$$

$$\Rightarrow \alpha p(\alpha x|I) = p(x|I)$$

$$f(x) = \alpha f(\alpha x)$$

$$f(x) = \frac{C}{x}$$

\equiv Scale Invariant Jeffreys Prior for Scale Parameters ($1/x$ prior)

- Limits x_{\min} and x_{\max} needed to normalise.

$$\int_{x_{\min}}^{x_{\max}} \frac{C}{x} dx = 1 = C \log \left(\frac{x_{\max}}{x_{\min}} \right)$$

- Alternatively suppose we are measuring location of something

$p(R|I)$ = position is at radius r to radi from sun

- Change origin $R' = R + C$:

$$p(R|I) dR = p(R'|I) dR'$$

$$p(R|I) = p(R + C|I)$$

- Then with limits x_{\min} and x_{\max} :

$$\int_{x_{\min}}^{x_{\max}} C = C(x_{\max} - x_{\min}) = 1$$

$$C = \frac{1}{x_{\max} - x_{\min}}$$

$$p(R) = \text{constant}$$

- Location invariant prior – uniform prior.

20.5 Nuisance Parameters and Marginalization

- Start:

$$P(\theta_1, \theta_2, \dots, \theta_n | D, I)$$

$$y = mx + b$$

$$P(mb | DI)$$

want:

$$P(m | DI) = \int P(m, b | DI) db$$

- Note that:

$$P(x, y) \Rightarrow P(x) = \int P(x, y) dy$$

- Marginalization is integrating out nuisance parameters.

- \hat{m}, \hat{b}

$$P(\theta_1, \theta_2 | D, I) = p(D | \theta_1, \theta_2, I) p(\theta_1, \theta_2) / \left(\int (\cdot) d\theta_1 d\theta_2 \right)$$

- Sample of θ_1 , and θ_2 , then we can marginalize numerically by ignoring θ_2 values.
- In this case we get a posterior distribution for θ_1 .
- Lots of ways to sample:
 - Grid Sampling
 - Random Sampling
 - MCMC Sampling (Markov Chain Monte Carlo)
- We will continue this next lecture.
- Then we will have all of the machinery to solve problems in any way we want:
- Bayesian, LS, MLE, MAP, MaxEnt, Frequentist, etc.

21 Tuesday November 25th 2025

21.1 Inverse Transform Sampling

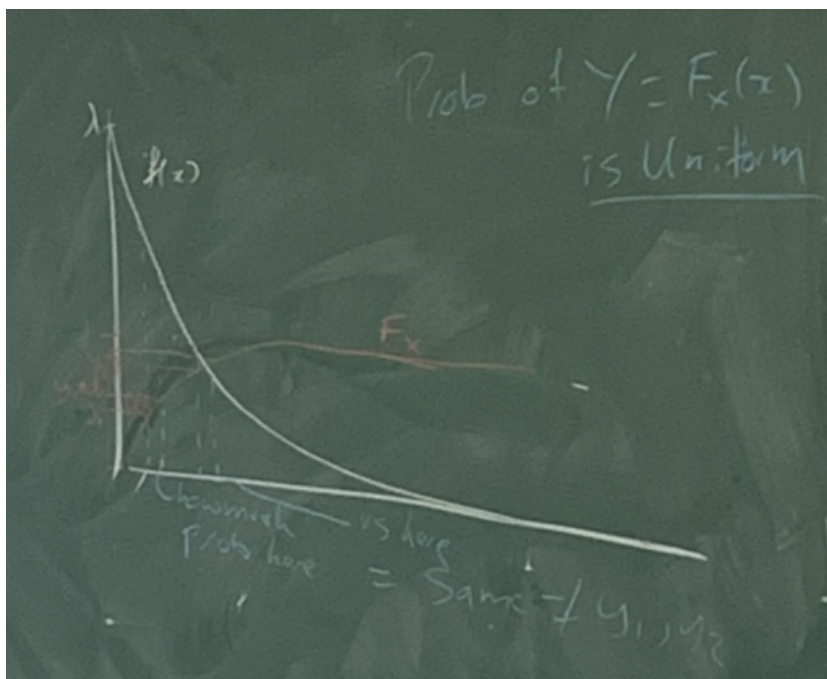
- We want to generate a sampling from a given pdf
- If available it is usually the best (but not usually possible...)
- Idea: Consider the CDF F_x of some PDF $f(x)$, $x \sim f(x)$, e.g. $f(x) = \lambda e^{-\lambda x}$

$$F_x(x) = \int_{-\infty}^x f(t) dt$$

$$= \Pr(X \leq x)$$

$$F_x : [0, 1]$$

- Let $y = F_x(x)$ random variable:
- choose x by $f(x)$, then get $y = F_x x$



- Prob of $y = F_x(x)$ is uniform
- Inverse F_x^{-1} will map y to x
- $\Pr(Y) \sim U(0, 1)$
- Uniform distribution satisfied
- $f_y(y) = \frac{dF_y(y)}{dy}$
- What is $F_y(y)$? $y \sim U(0, 1)$

$$\begin{aligned} F_y(y) &= \Pr(Y \leq y) \\ &= \Pr(F_x(x) \leq y) \end{aligned}$$

$$\begin{aligned}
&= \Pr(X \leq F_x^{-1}(y)) \\
&= F_x(F_x^{-1}(y)) = y
\end{aligned}$$

- F is monotonic increasing
- Let $Y = F_x^{-1}(x)$

$$X = F'(Y), \quad \text{where } Y \sim U(0,1)$$

- Pick u from $U(0,1)$
- Calculate $x = F_x^{-1}(u)$
- this will be distributed as $X \sim f(x)$

$$f(x) = \lambda e^{-\lambda x}$$

- CDF:

$$\begin{aligned}
F_x(x) &= \int_0^x \lambda e^{-\lambda t} dt \\
&= 1 - e^{-\lambda x}
\end{aligned}$$

- Pick $u \sim U(0,1)$

$$\begin{aligned}
u &= F_x(x) = 1 - e^{-\lambda x} \\
e^{-\lambda x} &= 1 - u \\
-\lambda x &= \log(1 - u)
\end{aligned}$$

$$x = -\frac{1}{\lambda} \log(1 - u)$$

21.2 Example: Sine Distribution

- Consider:

$$f(x) = \frac{1}{2} \sin x \quad x \in [0, \pi]$$

- CDF:

$$\begin{aligned}
F_x(x) &= \int_0^x \frac{1}{2} \sin t dt \\
&= \frac{1}{2}(1 - \cos x)
\end{aligned}$$

$$F(F^{-1}(u)) = u$$

$$u = \frac{1}{2}(1 - \cos F^{-1}(u))$$

$$\cos F^{-1}(u) = 1 - 2u$$

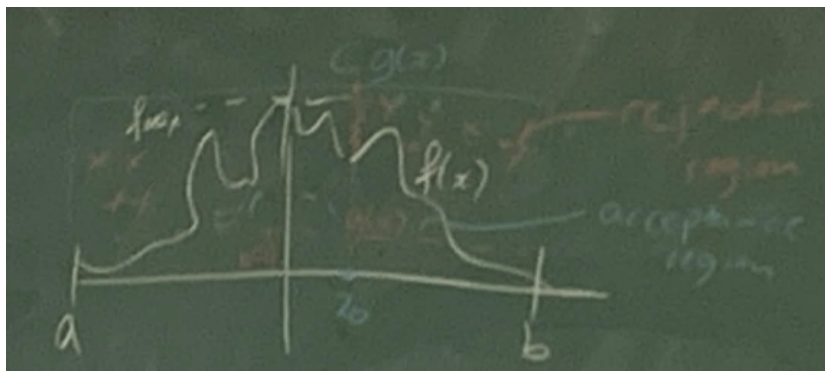
$$F^{-1}(u) = \arccos(1 - 2u)$$

21.3 Rejection Sampling (Accept-Reject Method)

- Usually distribution you want to sample from is complicated

21.4 Basic Rejection Sampling with Uniform Proposal

- Attempt 1: Make a sample over a finite domain of your variable
- Want process to generate samples $\sim f(x)$ $x \in [a, b]$



1. Generate x_0 according to $u(a, b)$:

$$g(x) = \frac{1}{b-a} = u(a, b) \equiv \text{proposal function}$$

2. Consider the function:

$$Cg(x) \text{ where } C \text{ is such that } Cg(x) \geq f(x) \quad \forall x \in [a, b]$$

For our case $C = (b-a)f_{\max}$

$$C = (b-a)f_{\max}$$

$$Cg(x) = f_{\max}$$

3. Generate a uniform number from $[0, Cg(x_0)]$, $U(0, Cg(x_0))$

4. If $u \leq f(x_0)$, accept x_0 If $u > f(x_0)$, reject x_0 .

$$U(0, Cg(x_0)) = Cg(x)U(0, 1)$$

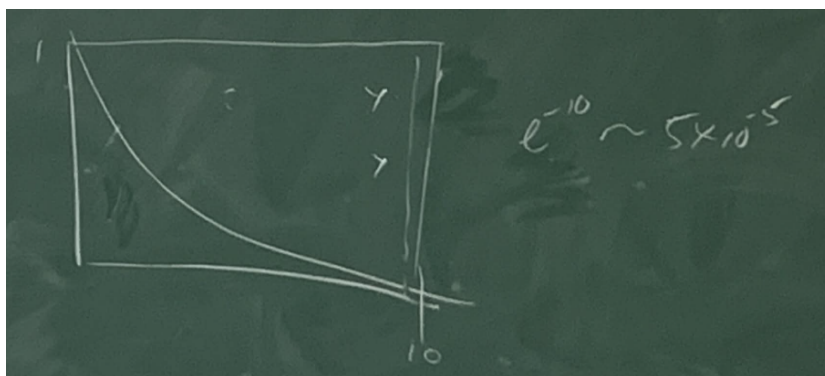
$$u' \sim U(0, 1)$$

$$u' = Cg(x_0) \leq f(x_0)$$

$$u' \leq \frac{f(x_0)}{Cg(x_0)}$$

Any C s.t. $Cg(x) \geq f(x)$ works, but you waste computing by picking C too large.

- Imagine this for e^{-x} on $[0, 10]$



21.5 General Rejection Sampling with Arbitrary Proposal

- Attempt 2: Reduce waste, modify proposal function

1. We can sample from it easily
2. It is closer in shape to f
3. $\exists C$ such that $Cg(x) \geq f(x) \forall x$

1. Sample $x_0 \sim g(x)$ proposal
2. throw uniform $u \sim U(0, 1)$
3. Accept x_0 based upon:

– Accept:

$$u \leq \frac{f(x_0)}{Cg(x_0)}$$

– Reject:

$$u > \frac{f(x_0)}{Cg(x_0)}$$

- A = accept point x selected from $g(x)$ with probability:

$$p = \frac{f(x)}{Cg(x)}$$

- Probability density of x given accept:

$$p(x|A) = \frac{p(A|x)g(x)}{p(A)} = \frac{\frac{f(x)}{Cg(x)}g(x)}{\int_{-\infty}^{\infty} \frac{f(x)}{Cg(x)}g(x) dx} = \frac{f(x)}{\int_{-\infty}^{\infty} f(x) dx}$$

22 Thursday, November 27th 2024

22.1 Curse of Dimensionality

- curse of dimensionality in this context means that a uniform distribution in high dimension with finite points will mostly be empty space on the inside, most will be near the boundary.

22.2 Monte Carlo Integration

- Monte Carlo Integration: often we need:

$$\int f(\vec{x}) d\vec{x}$$

especially hard in high dimensions (\mathbb{R}^n)

$$\int_a^b f(x) dx = \lim_{\Delta x \rightarrow 0} \sum f(x_i^*) \Delta x$$

where x_i^* is x in bin i .

- Expectation of f over pdf $p(x)$:

$$E(f) = \int f(x) p(x) dx$$

- Take:

$$p = U(a, b) = \frac{1}{b - a}$$

then:

$$E_a(f) = \int_a^b f(x) \frac{1}{b-a} dx = \frac{1}{b-a} \int_a^b f(x) dx$$

so:

$$\int_a^b f(x) dx = (b-a) E_{U(a,b)}(f)$$

- In general if $x \sim X$
- If we take sample of size n

$$\hat{\mu} = \text{sample mean} = \frac{1}{n} \sum_{i=1}^n x_i$$

is an unbiased consistent estimator of true sample mean

- As $n \rightarrow \infty$,

$$\lim \frac{1}{n} \sum_{i=1}^n x_i \rightarrow \mu$$

- Also true for functions of x

$$E_x(f) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n f(x_i)$$

where $x_i \sim X$

- So back to our equation from before we have:

$$\int_a^b f(x) dx = (b-a) E_{U(a,b)}(f) = (b-a) \frac{1}{n} \lim_{n \rightarrow \infty} \sum_{i=1}^n f(x_i)$$

where $x_i \sim U(a, b)$

- c.f. (cumulative frequency?)

$$\left(\frac{b-a}{n}\right) \sum_i f(x_i)$$

where the prefactor is the width of each bin (Δx)

- Many dimensions: Riemann: 1000 bins, n-dims $f(\vec{x})$, 10^{3n} bins total, so even for a small $n=3$, that gives you 10^9 bins.

22.3 Importance Sampling

- Importance Sampling:

$$I = \int f(x) dx = \int f(x) \frac{q(x)}{q(x)} dx = \int \frac{f(x)}{q(x)} q(x) dx$$

- Trick: $q(x)$ is probability density with $q \neq 0$ (at least where $f \neq 0$)

$$I_n = E_q \left(\frac{f(x)}{q(x)} \right)$$

$$\approx \frac{1}{n} \sum_{i=1}^n \frac{f(x_i)}{q(x_i)}$$

where $x_i \sim q(x)$

$$I = \frac{1}{n} \sum_{i=1}^n w(x_i)$$

where $w(x) = \frac{f(x)}{q(x)}$

$$E(\hat{I}_n) = \frac{1}{n} \sum \left(\frac{f(x_i)}{q(x_i)} \right) = \frac{1}{n} = I$$

- Variance of \hat{I}_n :
- General $V(y_1 + y_2 + \dots + y_n) = V(y_1) + V(y_2) + \dots + V(y_n)$ if y_i independent

$$\begin{aligned} V_q(\hat{I}_n) &= \frac{1}{n^2} V_q \left(\sum_{i=1}^n w(x_i) \right) = \frac{1}{n^2} n V_q(w(x)) = \frac{1}{n} V_q(w(x)) \\ &= \frac{1}{n} \left(E_q(w(x)^2) - E_q(w(x))^2 \right) \\ &= \frac{1}{n} \left(E_q(w(x)^2) - I^2 \right) \\ &= \frac{1}{n} \left(\int \frac{f(x)^2}{q(x)} dx - I^2 \right) \end{aligned}$$

- Crude Monte Carlo vs Importance Sampling
 - Crude MC: sample from uniform distribution over domain of integration
 - IS: sample from distribution that mimics behavior of $f(x)$
- Minimize this – q to be proportional to $|f(x)|$
- Constraint is that $q(x)$ is a valid pdf: $\int q(x) dx = 1$
- f large, q large, f small, q small.

22.4 Importance Sampling for Bayesian Inference

- Moving on:

$$\int p(D|\theta)p(\theta)d\theta$$

norm in Bayes

$$p(\theta|D) = \frac{p(D|\theta)p(\theta)}{\int p(D|\theta)p(\theta)d\theta}$$

Need:

$$E_{p(\theta|D)}(f) = \int f(\theta)p(\theta|D)d\theta$$

- Also:

$$\int \sigma(m)p(m|D)dm$$

- The integral before $\int f(x)p(x)dx$ is called an expectation value of f with respect to p $E_p(f)$

$$E_p(f) = \int f(x)p(x)dx = \int f(x)\frac{p(x)}{q(x)}q(x)dx = E_q\left(f(x)\frac{p(x)}{q(x)}\right)$$

$$E_p(f) = E_q\left(f(x)\frac{p(x)}{q(x)}\right)$$

$$E_p(f) \approx \frac{1}{n} \sum_{i=1}^n f(x_i^q) \frac{p(x_i^q)}{q(x_i^q)}$$

$$= \frac{1}{n} \sum_{i=1}^n f(x_i^q)w(x_i^q)$$

where $x_i^q \sim q(x)$ and $w = \frac{p}{q}$

- suppose p right shape but $\int p(x)dx = Z \neq 1$

$$\tilde{p}(x) = \frac{p(x)}{Z} \quad \text{is a valid pdf}$$

$$E_{\tilde{p}}(f) = \int f(x)\tilde{p}(x)dx$$

$$= E_q\left(f(x)\frac{\tilde{p}(x)}{q(x)}\right)$$

$$= \frac{1}{Z} E_q\left(f(x)\frac{p(x)}{q(x)}\right)$$

$$= \frac{1}{Z} E_q(f(x)w(x))$$

- \tilde{p} is a pdf

$$E_{\tilde{p}}(1) = 1 = \frac{1}{Z} E_q(w(x))$$

- So:

$$Z = E_q(w(x))$$

$$E_{\tilde{p}}(f) = \frac{E_q(f(x)w(x))}{E_q(w(x))} \approx \frac{\frac{1}{n} \sum_{i=1}^n f(x_i^q) w(x_i^q)}{\frac{1}{n} \sum_{i=1}^n w(x_i^q)}$$

where $x_i^q \sim q(x)$

$$= \frac{\bar{w}f}{\bar{w}}$$

where $\bar{\cdot}$ is the average.

22.5 Metropolis-Hastings Algorithm

- Metropolis (Hastings) Algorithm
- Use to sample from any (usually complicated) pdf $p(x)$
- Don't need normalized pdf!
- Proposal function q ,

$$p(\theta|\vec{x}) = p(\vec{x}|\theta)p(\theta)$$

- But it will depend on last accepted proposal
- Algorithm:
 - Start with point \vec{x}_0
 - Pick \vec{x}' from proposal distribution $q(\vec{x}'|\vec{x}_0)$
 - Metropolis q is symmetric: $q(\vec{a}|\vec{b}) = q(\vec{b}|\vec{a})$
- Example: multi dimensional Gaussian:

$$\begin{aligned} q(\vec{a}|\vec{b}) &= \vec{b} + N(0, \Sigma) \\ &= N(\vec{b}, \Sigma) \end{aligned}$$

- Hastings q can be not symmetric.
- Steps:
 1. Start with point \vec{x}_0
 2. Generate proposal $\vec{x}' \sim q(\vec{x}'|\vec{x}_0)$
 3. Form acceptance ratio:

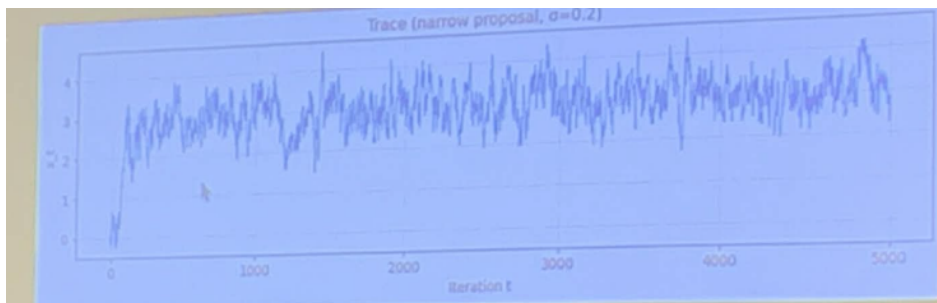
$$r = \frac{p(\vec{x}')q(\vec{x}_0|\vec{x}')}{p(\vec{x}_0)q(\vec{x}'|\vec{x}_0)}$$


```

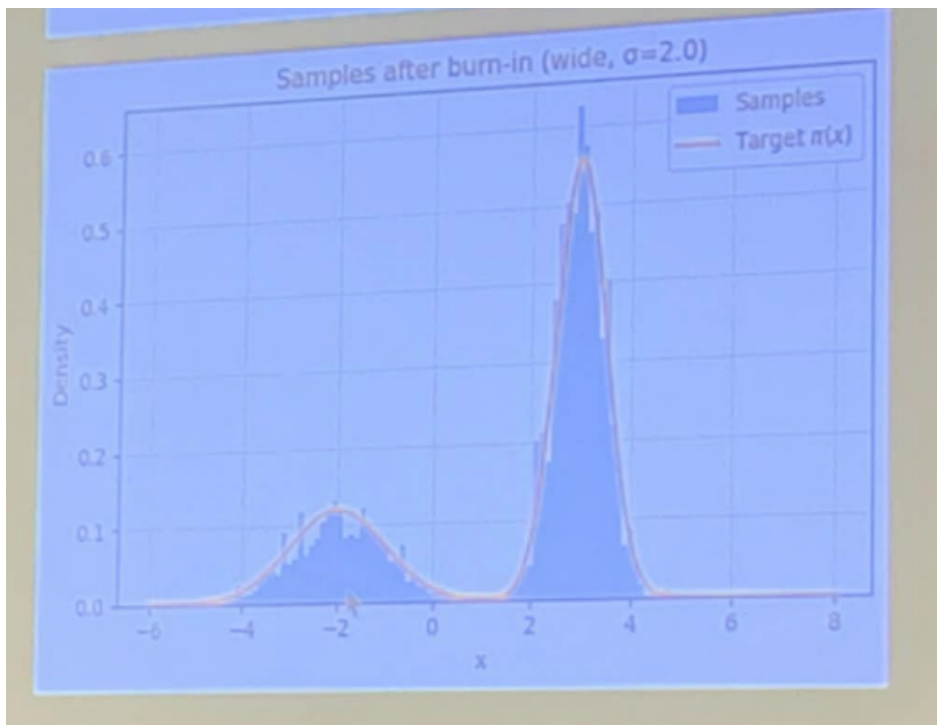
4.      If  $r \geq 1$ :
          accept proposal:  $\vec{x}_{t+1} = \vec{x}_t'$ 
      If  $r < 1$ :
          accept with probability  $r$  (throw unit number  $x_{t+1} = x_t$ )
      else:
           $x_{t+1} = x_t$  (stay at same place)

```

- Algorithm is a stochastic process called Markov Process,
- Output: $\{x_0, x_1, x_2, \dots, x_n\} = \text{Markov Chain}$
- Output (after burn-in steps) is samples from $\sim p(x)$



- Here we can see the burn-in period at the start



- MCMC for a complicated bi-modal distribution