

# IEDA 2540 Statistics for Engineers

## Introduction

Wei YOU



香港科技大學

THE HONG KONG UNIVERSITY OF  
SCIENCE AND TECHNOLOGY

Spring, 2025

# What is Statistics?

The thinking involved in statistics is not new!

- Imagine meeting two friends: one is 1.6m tall, the other 1.8m tall.
- Based on your daily experience, you'd likely guess the taller one is a man.
- Could you be wrong? Sure!
- Yet, you wouldn't hesitate to trust your gut.

# Statistical Way of Thinking

- Statistics **summarize** past experience.
- When facing a new scenario, we **generalize/predict** from previous experience by making an informed guess.

### Example:

- **Summarize:** “On *average*, I cycle about 100 miles a week.”
- **Predict:** “We can *expect* a lot of rain at this time of year.”
- **Predict:** “The earlier you start preparing, the *more likely* you are to do well in exams.”

## Statistical Thinking in Everyday Life

Humans learn from the environment to gain greater control.

- However, we can never know everything for certainty.
- Fortunately, we *observe* and notice *patterns* and *regularities*.

**Example:** Counting and measuring in a *rough-and-ready* fashion.

- Humans have the intuition to rely on habits of “quantification”: “How often?”, “how big?”, “how difficult?”, “how far?”.

Summarizing one aspect of one subject.

## Statistical Thinking in Everyday Life

**Example:** Noticing several aspects of the same subject

- the size of a potato-crop in a particular field this year
- Average temperature of the growing season.
- Sunshine.
- Amount of fertilizer used.
- Nature of soil (acidity or basicity).
- Rainfall.

**Example:** More interestingly, we compare several subjects with lots in common but differing in some respects.

- How the size of a potato varies between several different fields in the same year?
- Summarizing multiple aspects of multiple subjects → **data collection**.

## Looking for Patterns

By instinct, when facing a collection of observations, we look for **connections**, **patterns**, **similarities** and **differences**.

**Example:** Studying the size of potato-crop harvested from different fields

- Are there differences in their sizes?
  - If one potato from field A is larger than one potato from field B, can we say field A is better?
  - If the average of 5 potatoes from field A is larger than that from field B, is field A better?
  - How confident are we in such claims?
- What are the differences between the fields? What could be the cause of the size difference?
  - Soil? Weather? Planting methods? Or a combination of several factors?

## From Patterns to Action

More vital questions

- What can we learn from the patterns we identified?
- How do they help us act more efficiently in the future?

**This is where statistics comes in.**

- **Statistics** is a scientific way to:
  - **collect** and **summarize** observations,
  - **identify** patterns,
  - **generalize** conclusions,

and avoid jumping to **hasty conclusions** by considering what could have happen or what could have gone wrong.

**To which extent can we generalize from our limited experiences?**

## Be Cautious about Generalizing from Experience

Tendency to generalize is human nature.

**Example:** When we generously apply fertilizer, we harvest bigger than usual potato.

- Shall we apply more fertilizer?
- Shall we do the same for other plants?
- Shall we do the same for other fields?

Would you think it safe to generalize in this way – on the basis of experience with one field? Why, or why not?

- What is true for one subject may not be true for others!
- To generalize more confidently, we need **more experience/data!**



## Dealing with Chances

Even if we have collected as much data as we can, is it possible to predict with 100% certainty?

Think about the following two arguments. Which do you think is more likely to be correct?

- “I predict *most* fields will have bigger potatoes if we do such-and-such.”
- “I predict that *any specific* field will have bigger potatoes if we do such-and-such.”

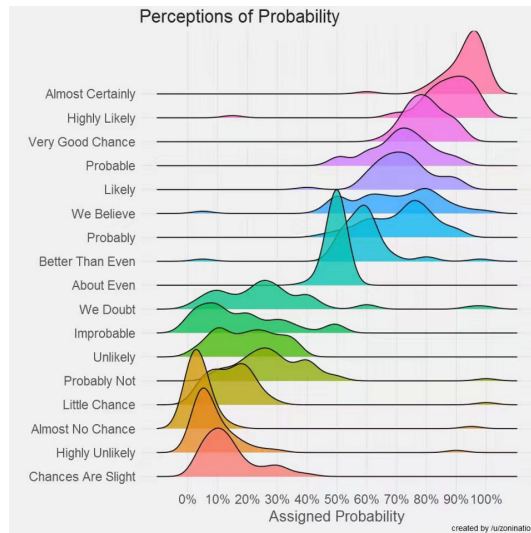
There is no 100% certainty, especially when dealing with people, things and events!

- Statistics helps us to look for **reliable regularities** and associations among things “in general,” “in the long run” and “with high chance.”
- At the same time, it teaches us to **manage expectation** about our predictions, especially when it comes to specific individuals.

# What is Probability?

Probability is the language of uncertainty and variability.

- Probability used in daily communication: “Almost certainly”, “probably not”, ...
- What are the meaning of these phrases?



# What is Probability?

**Example:** Probability used in daily communication

- When some one says “highly likely,” what does it mean?
- With a probability approach, we need to assume a “distribution,” i.e., how likely a certain meaning of “highly likely” occurs.

**Example:** Imagine a colleague is always late for meetings. Someone says, “He’ll be on time for once, *highly likely*,” sarcastically.

- What is the meaning of “highly likely” out of an average person’s mouth? Mean.
- How variable is the meaning of “highly likely?” Variance.

## History of Probability

A 17th century gambler Chevalier de Méré

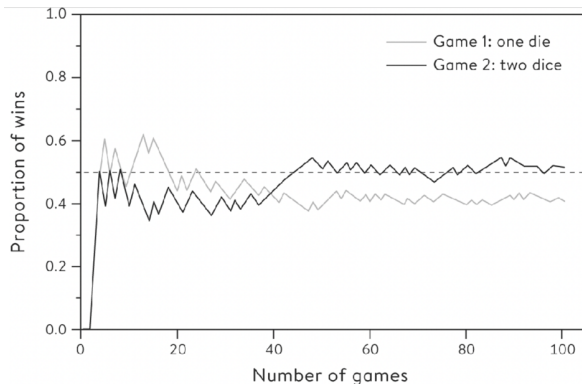
- **Game 1:** "Throw a fair die at most 4 times and win if you get a six."
- **Game 2:** "Throw two fair dice at most 24 times and win if you get a double-six."

Which one would you bet on?



## Chevalier de Méré's Reasoning

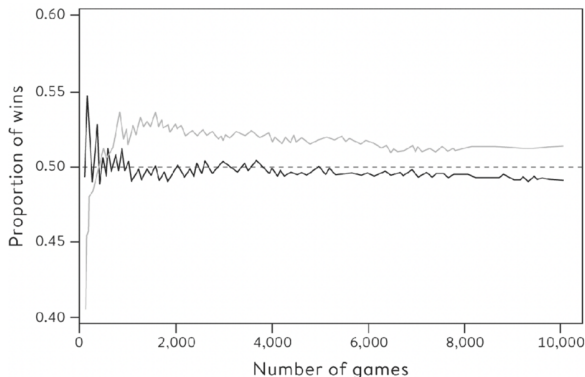
He simulated the game 100 times, and found that Game 2 has higher reward.



## Chevalier de Méré's Reasoning

What went wrong?

- In the presence of randomness, we need to wisely interpret the experiments we did.



## The Math Behind

Let  $p = \frac{1}{6}$  be the probability of a six in one roll.

- **Game 1:** “Throw a fair die at most 4 times and win if you get a six.”
  - Probability of winning:

$$1 - (1 - p)^4 = 0.518.$$

- **Game 2:** “Throw two fair dice at most 24 times and win if you get a double-six.”
  - Probability of winning:

$$1 - \left(1 - \frac{p}{6}\right)^{4 \times 6} = 1 - \left[\left(1 - \frac{p}{6}\right)^6\right]^4 = 0.491.$$

Read more here: <https://tinyurl.com/3t3y9xku>

## History of Probability: Key Figures

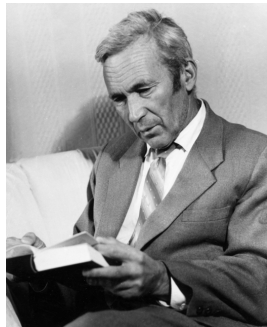
- Chevalier de Méré complained to **Pascal** and **Fermat**, who laid the foundation of probability theory.
- Probability theory developed significantly through the 17th to 19th centuries.
- **Andrey Kolmogorov** is considered the founder of modern probability theory.



Blaise Pascal



Pierre de Fermat





## History of Statistics

Statistics, as a word, is derived from the German word “Statistik” meaning “description of a state/country.”

- Systematic collection of data on population and economy dates back to Renaissance (14th - 16th century).
- John Graunt (British, founder of demography, epidemiologist), natural and political observations made upon the **bills of mortality**, 1662.
- More and more regions started to systematically collect data from 1700s.
- Statistics became concerned with **inferencing conclusions** from a sample of numerical data, late 1800s.

### Example: Bill of mortality

TABLE 1.2 *John Graunt's Mortality Table*

Age at Death	Number of Deaths per 100 Births
0-6	36
6-16	24
16-26	15
26-36	9
36-46	6
46-56	4
56-66	3
66-76	2
76 and greater	1

*Note: The categories go up to but do not include the right-hand value. For instance, 0-6 means all ages from 0 up through 5.*

The Table of CASUALTIES.

The Year of our Lord	1647	1648	1649	1650	1651	1652	1653	1654	1655	1656	1657	1658	1659	1660	1661	1662	1663	1664	1665	1666	1667	1668	1669	1670	1671	1672	1673	1674	1675	1676	1677	1678	1679	1680	1681	1682	1683	1684	1685	1686	1687	1688	1689	1690	1691	1692	1693	1694	1695	1696	1697	1698	1699	1700	1701	1702	1703	1704	1705	1706	1707	1708	1709	1710	1711	1712	1713	1714	1715	1716	1717	1718	1719	1720	1721	1722	1723	1724	1725	1726	1727	1728	1729	1730	1731	1732	1733	1734	1735	1736	1737	1738	1739	1740	1741	1742	1743	1744	1745	1746	1747	1748	1749	1750	1751	1752	1753	1754	1755	1756	1757	1758	1759	1760	1761	1762	1763	1764	1765	1766	1767	1768	1769	1770	1771	1772	1773	1774	1775	1776	1777	1778	1779	1780	1781	1782	1783	1784	1785	1786	1787	1788	1789	1790	1791	1792	1793	1794	1795	1796	1797	1798	1799	1800	1801	1802	1803	1804	1805	1806	1807	1808	1809	1810	1811	1812	1813	1814	1815	1816	1817	1818	1819	1820	1821	1822	1823	1824	1825	1826	1827	1828	1829	1830	1831	1832	1833	1834	1835	1836	1837	1838	1839	1840	1841	1842	1843	1844	1845	1846	1847	1848	1849	1850	1851	1852	1853	1854	1855	1856	1857	1858	1859	1860	1861	1862	1863	1864	1865	1866	1867	1868	1869	1870	1871	1872	1873	1874	1875	1876	1877	1878	1879	1880	1881	1882	1883	1884	1885	1886	1887	1888	1889	1890	1891	1892	1893	1894	1895	1896	1897	1898	1899	1900	1901	1902	1903	1904	1905	1906	1907	1908	1909	1910	1911	1912	1913	1914	1915	1916	1917	1918	1919	1920	1921	1922	1923	1924	1925	1926	1927	1928	1929	1930	1931	1932	1933	1934	1935	1936	1937	1938	1939	1940	1941	1942	1943	1944	1945	1946	1947	1948	1949	1950	1951	1952	1953	1954	1955	1956	1957	1958	1959	1960	1961	1962	1963	1964	1965	1966	1967	1968	1969	1970	1971	1972	1973	1974	1975	1976	1977	1978	1979	1980	1981	1982	1983	1984	1985	1986	1987	1988	1989	1990	1991	1992	1993	1994	1995	1996	1997	1998	1999	2000	2001	2002	2003	2004	2005	2006	2007	2008	2009	2010	2011	2012	2013	2014	2015	2016	2017	2018	2019	2020	2021	2022	2023	2024	2025	2026	2027	2028	2029	2030	2031	2032	2033	2034	2035	2036	2037	2038	2039	2040	2041	2042	2043	2044	2045	2046	2047	2048	2049	2050	2051	2052	2053	2054	2055	2056	2057	2058	2059	2060	2061	2062	2063	2064	2065	2066	2067	2068	2069	2070	2071	2072	2073	2074	2075	2076	2077	2078	2079	2080	2081	2082	2083	2084	2085	2086	2087	2088	2089	2090	2091	2092	2093	2094	2095	2096	2097	2098	2099	2100	2101	2102	2103	2104	2105	2106	2107	2108	2109	2110	2111	2112	2113	2114	2115	2116	2117	2118	2119	2120	2121	2122	2123	2124	2125	2126	2127	2128	2129	2130	2131	2132	2133	2134	2135	2136	2137	2138	2139	2140	2141	2142	2143	2144	2145	2146	2147	2148	2149	2150	2151	2152	2153	2154	2155	2156	2157	2158	2159	2160	2161	2162	2163	2164	2165	2166	2167	2168	2169	2170	2171	2172	2173	2174	2175	2176	2177	2178	2179	2180	2181	2182	2183	2184	2185	2186	2187	2188	2189	2190	2191	2192	2193	2194	2195	2196	2197	2198	2199	2200	2201	2202	2203	2204	2205	2206	2207	2208	2209	2210	2211	2212	2213	2214	2215	2216	2217	2218	2219	2220	2221	2222	2223	2224	2225	2226	2227	2228	2229	2230	2231	2232	2233	2234	2235	2236	2237	2238	2239	2240	2241	2242	2243	2244	2245	2246	2247	2248	2249	2250	2251	2252	2253	2254	2255	2256	2257	2258	2259	2260	2261	2262	2263	2264	2265	2266	2267	2268	2269	2270	2271	2272	2273	2274	2275	2276	2277	2278	2279	2280	2281	2282	2283	2284	2285	2286	2287	2288	2289	2290	2291	2292	2293	2294	2295	2296	2297	2298	2299	2300	2301	2302	2303	2304	2305	2306	2307	2308	2309	2310	2311	2312	2313	2314	2315	2316	2317	2318	2319	2320	2321	2322	2323	2324	2325	2326	2327	2328	2329	2330	2331	2332	2333	2334	2335	2336	2337	2338	2339	2340	2341	2342	2343	2344	2345	2346	2347	2348	2349	2350	2351	2352	2353	2354	2355	2356	2357	2358	2359	2360	2361	2362	2363	2364	2365	2366	2367	2368	2369	2370	2371	2372	2373	2374	2375	2376	2377	2378	2379	2380	2381	2382	2383	2384	2385	2386	2387	2388	2389	2390	2391	2392	2393	2394	2395	2396	2397	2398	2399	2400	2401	2402	2403	2404	2405	2406	2407	2408	2409	2410	2411	2412	2413	2414	2415	2416	2417	2418	2419	2420	2421	2422	2423	2424	2425	2426	2427	2428	2429	2430	2431	2432	2433	2434	2435	2436	2437	2438	2439	2440	2441	2442	2443	2444	2445	2446	2447	2448	2449	2450	2451	2452	2453	2454	2455	2456	2457	2458	2459	2460	2461	2462	2463	2464	2465	2466	2467	2468	2469	2470	2471	2472	2473	2474	2475	2476	2477	2478	2479	2480	2481	2482	2483	2484	2485	2486	2487	2488	2489	2490	2491	2492	2493	2494	2495	2496	2497	2498	2499	2500	2501	2502	2503	2504	2505	2506	2507	2508	2509	2510	2511	2512	2513	2514	2515	2516	2517	2518	2519	2520	2521	2522	2523	2524	2525	2526	2527	2528	2529	2530	2531	2532	2533	2534	2535	2536	2537	2538	2539	2540	2541	2542	2543	2544	2545	2546	2547	2548	2549	2550	2551	2552	2553	2554	2555	2556	2557	2558	2559	2560	2561	2562	2563	2564	2565	2566	2567	2568	2569	2570	2571	2572	2573	2574	2575	2576	2577	2578	2579	2580	2581	2582	2583	2584	2585	2586	2587	2588	2589	2590	2591	2592	2593	2594	2595	2596	2597	2598	2599	2600	2601	2602	2603	2604	2605	2606	2607	2608	2609	2610	2611	2612	2613	2614	2615	2616	2617	2618	2619	2620	2621	2622	2623	2624	2625	2626	2627	2628	2629	2630	2631	2632	2633	2634	2635	2636	2637	2638	2639	2640	2641	2642	2643	2644	2645	2646	2647	2648	2649	2650	2651	2652	2653	2654	2655	2656	2657	2658	2659	2660	2661	2662	2663	2664	2665	2666	2667	2668	2669	2670	2671	2672	2673	2674	2675	2676	2677	2678	2679	2680	2681	2682	2683	2684	2685	2686	2687	2688	2689	2690	2691	2692	2693	2694	2695	2696	2697	2698	2699	2700	2701	2702	2703	2704	2705	2706	2707	2708	2709	2710	2711	2712	2713	2714	2715	2716	2717	2718	2719	2720	2721	2722	2723	2724	2725	2726	2727	2728	2729	2730	2731	2732	2733	2734	2735	2736	2737	2738	2739	2740	2741	2742	2743	2744	2745	2746	2747	2748	2749	2750	2751	2752	2753	2754	2755	2756	2757	2758	2759	2760	2761	2762	2763	2764	2765	2766	2767	2768	2769	2770	2771	2772	2773	2774	2775	2776	2777	2778	2779	2780	2781	2782	2783	2784	2785	2786	2787	2788	2789	2790	2791	2792	2793	2794	2795	2796	2797	2798	2799	2800	2801	2802	2803	2804	2805	2806	2807	2808	2809	2810	2811	2812	2813	2814	2815	2816	2817	2818	2819	2820	2821	2822	2823	2824	2825	2826	2827	2828	2829	2830	2831	2832	2833	2834	2835	2836	2837	2838	2839	2840	2841	2842	2843	2844	2845	2846	2847	2848	2849	2850	2851	2852	2853	2854	2855	2856	2857	2858	2859	2860	2861	2862	2863	2864	2865	2866	2867	2868	2869	2870	2871	2872	2873	2874	2875	2876	2877	2878	2879	2880	2881	2882	2883	2884	2885	2886	2887	2888	2889	2890	2891	2892	2893	2894	2895	2896	2897	2898	2899	2900	2901	2902	2903	2904	2905	2906	2907	2908	2909	2910	2911	2912	2913	2914	2915	2916	2917	2918	2919	2920	2921	2922	2923	2924	2925	2926	2927	2928	2929	2930	2931	2932	2933	2934	2935	2936	2937	2938	2939	2940	2941	2942	2943	2944	2945	2946	2947	2948	2949	2950	2951	2952	2953	2954	2955	2956	2957	2958	2959	2960	2961	2962	2963	2964	2965	2966	2967	2968	2969	2970	2971	2972	2973	2974	2975	2976	2977	2978	2979	2980	2981	2982	2983	2984	2985	2986	2987	2988	2989	2990	2991	2992	2993	2994	2995	2996	2997	2998	2999	3000
----------------------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------

Place the Table at page 114.

## Probability and Statistics Today: Virus Testing

A patient takes a virus test that returns positive. Consider:

- 0.3% of the population have Covid.
- PCR test correctly detects Covid with 95% chance.
- PCR test incorrectly detects Covid with 5% chance when absent.

What is the probability they really have the virus? Only about 5.4%!

<https://calculator.testingwisely.com/>

- How to design good tests? → Hypothesis testing.

The math:

$$\begin{aligned}\mathbb{P}(C|+) &= \frac{\mathbb{P}(+|C) \cdot \mathbb{P}(C)}{\mathbb{P}(+|C) \cdot \mathbb{P}(C) + \mathbb{P}(+|\neg C) \cdot \mathbb{P}(\neg C)} \\ &= \frac{0.95 \times 0.003}{0.95 \times 0.003 + 0.05 \times 0.997} \approx 0.054\end{aligned}$$

## Probability and Statistics Today: Online Retailers

Amazon uses statistics to understand and predict your preferences based on your browsing and purchasing history.

- **Data Collection:**

- Tracks items you view, buy, and review.
- Aggregates data from millions of users.

- **Pattern Identification:**

- Uses statistical methods to find similarities between users and products.
- Applies techniques like collaborative filtering and clustering.

- **Prediction:**

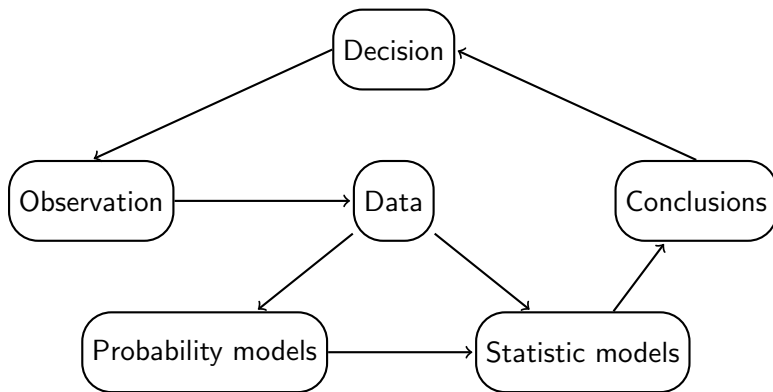
- Predicts products you might be interested in.
- Continuously updates recommendations as your behavior evolves.

## Netflix Prize (2006-2009)

Netflix launched a \$1 million contest to improve its movie recommendation system.

- **In-house algorithm:** Cinematch, a straightforward statistical linear model with extensive data conditioning (RMSE 0.9525).
- **Data context:** Before the era of “Big Data”, over 100 million ratings of 17,770 movies from 480,189 customers were available.
- **Winner:** Team “BellKor’s Pragmatic Chaos” achieved a 10% improvement with an RMSE of 0.8567.
- **Methods explored:** Regression, k-nearest-neighbor, matrix factorization, kernel ridge regression, neural networks, decision trees, and more.

## What is this course about?



Statistics enable us to summarize, compare and predict more precisely than we normally would in everyday conversation.