

Data Science & Big Data for Actuaries

Arthur Charpentier (Université de Rennes 1 & UQàM)

Universitat de Barcelona, April 2016.

<http://freakonometrics.hypotheses.org>



Data Science & Big Data for Actuaries

Arthur Charpentier (Université de Rennes 1 & UQàM)

Professor, Economics Department, Univ. Rennes 1

In charge of Data Science for Actuaries program, IA

Research Chair *actinfo* (Institut Louis Bachelier)

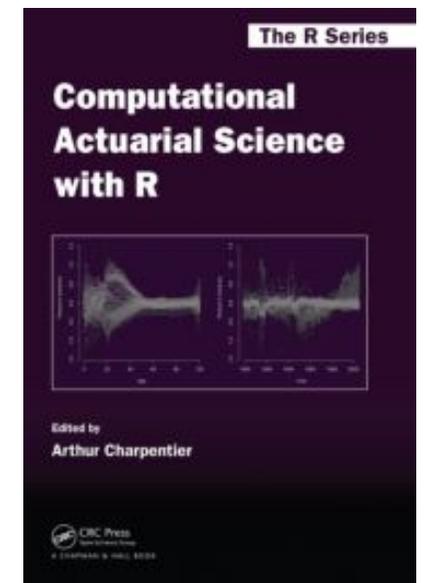
(previously Actuarial Sciences at UQàM & ENSAE Paristech
actuary in Hong Kong, IT & Stats FFSA)

PhD in Statistics (KU Leuven), Fellow Institute of Actuaries

MSc in Financial Mathematics (Paris Dauphine) & ENSAE

Editor of the freakonometrics.hypotheses.org's blog

Editor of Computational Actuarial Science, CRC



Data

“People use statistics as the drunken man uses lamp posts - for support rather than illumination”,
Andrew Lang or not

see also Chris Anderson *The End of Theory: The Data Deluge Makes the Scientific Method Obsolete*, 2008

1. **An Overview on (Big) Data**
2. **Big Data & Statistical/Machine Learning**
3. **Classification Models**
4. **Small Data & Bayesian Philosophy**
5. **Data, Models & Actuarial Science**



Part 1.

An Overview on (Big) Data



Historical Aspects of Data



Storing Data: Tally sticks, used starting in the Paleolithic area



A tally (or tally stick) was an ancient memory aid device used to record and document numbers, quantities, or even messages.

Historical Aspects of Data

The Table of CASUALTIES.

The Years 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	
-----------------------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	--

Historical Aspects of Data

Data Manipulation: Herman Hollerith created a Tabulating Machine that uses punch cards to reduce the workload of US Census, in 1881, see [1880 Census](#), $n = 50$ million Americans.

1	56 73 144	Fisk Fred	M 40		1	Labourer
2		— Emilia	W 30	Wife	1	Housewife
3		— Gustav	M 8	Son	1	At School
4		— Anna	W 4	Daughter	1	
5		— Richard	M 2	Son	1	
6	73 145	Hannichen Peter	M 38		1	Wagon Driver
7		— Ernestine	W 26	Wife	1	Housewife
8		— Emma	W 8	Daughter	1	At School
9		— Jakob	M 6	Son	1	At School
10		— Pauline	W 4	Daughter	1	
11		— Anna	W 2	Daughter	1	
12		— Willie	M 6	Son	1	

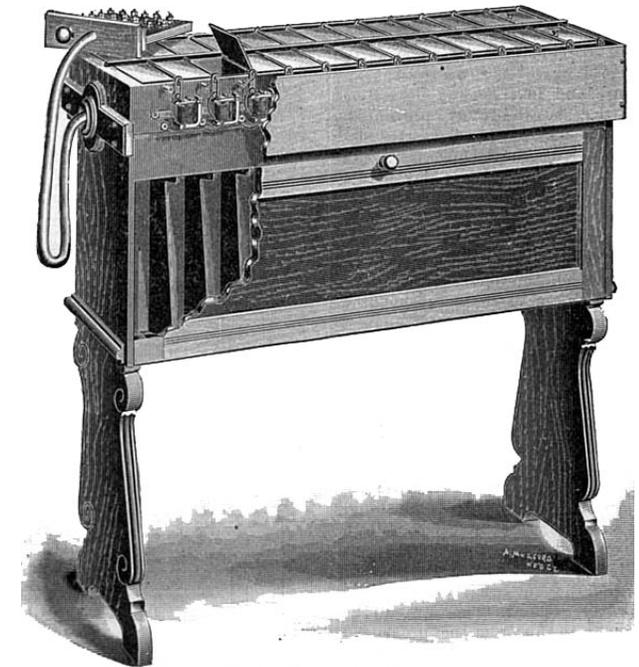


Fig. 3.—Sorting Machine.

Hollerith's Electric Sorting and Tabulating Machine.

Historical Aspects of Data

Survey and Polls: 1936 US elections

Literary Digest Poll based on 2.4 million readers

A. Landon: 57% vs. F.D. Roosevelt: 43%

George Gallup sample of about 50,000 people

A. Landon: 44% vs. F.D. Roosevelt: 56%

Actual results

A. Landon: 38% vs. F.D. Roosevelt: 62%

Sampling techniques, polls, predictions based on small samples



Historical Aspects of Data

Allen-Scott Report

Data Center Plan Called Privacy Invasion

By ROBERT S. ALLEN
and PAUL SCOTT

WASHINGTON — A special White House task force is recommending the creation of a federal data center which eventually could have a comprehensive file on every man, woman and child in the country.

Now under study in inner administration circles, the still-secret report advocates the gradual transfer of all governmental records and statistics to magnetic computer tape, which would be turned over to a newly-created agency that would function as a general data center.

The computerized information would be available, at the push of a button, to a wide range of government authorities.

Estimated cost of the pro-

cial Security, census data, medical, credit and criminal reports.

"Comprehensive information of this kind, centralized in one agency," says Gallagher, "could constitute a highly dangerous dossier bank. Such an agency would be a distinct departure from our American tradition."

Subcommittee investigators have ascertained that the task force's report states that a vast accumulation of government records already is on computer tape and could be turned over to the proposed general data center immediately. Listed as among these available files are:

Internal Revenue Service — 742 million personal and corporate tax returns.

Defense Department — 14

the most intimate information, the investigators learned, are freely passed around among agencies. Graphically illustrative of this practice and its harsh consequences are the following two instances:

A teenager visiting Washington stayed with an uncle, at his mother's suggestion. During the night the boy was sexually assaulted by the uncle. Years later, as a Phi Beta Kappa graduate from a leading Eastern university, the boy applied for a job with the National Security Agency. During a required lie detector test he told about the assault. His frank admission cost him the desired job.

But that wasn't all. This affair, in which he was an innocent victim, hounded him again

Data Center: The US Government plans the world's first data center to store 742 million tax returns and 175 million sets of fingerprints, in 1965.

Historical Aspects of Data



Historical Aspects of Data

Data Manipulation: Relational Database model developed by Edgar F. Codd

See [Relational Model of Data for Large Shared Data Banks](#), [Codd \(1970\)](#)

Considered as a major breakthrough for users and machine designers

Data or tables are thought as a matrix composed of intersecting rows and columns, each columns being attributes.

Tables are related to each other through a common attribute.

Concept of relational diagrams

The Two Cultures

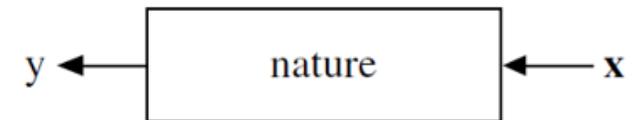
‘The Two Cultures’, see [Breiman \(2001\)](#)

- **Data Modeling** (statistics, econometrics)
- **Algorithmic Modeling** (computational & algorithmics)

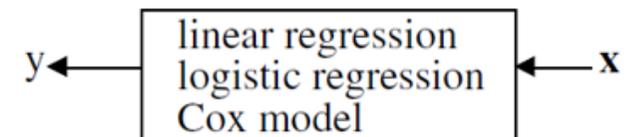
‘**Big Data Dynamic Factor Models for Macroeconomic Measurement and Forecasting**’, [Diebold \(2000\)](#)

Statistical Science
2001, Vol. 16, No. 3, 199–231

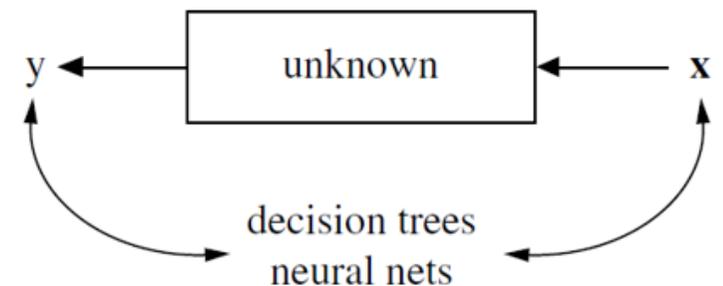
Statistical Modeling: The Two Cultures



The Data Modeling Culture



The Algorithmic Modeling Culture



And the XIXth Century...

Nature's special issue on Big Data, [Nature \(2008\)](#) and many of business journals



And the XIXth Century...

Technology changed, HDFS (Hadoop Distribution File System), MapReduce

The New York Times Technology | Personal Tech | Business Day Log In | Register Now

Bits

Search Bits

OCTOBER 24, 2012, 9:00 AM | 4 Comments

Big Data in More Hands

By QUENTIN HARDY

FACEBOOK

TWITTER

GOOGLE+

SAVE

E-MAIL

SHARE

PRINT

Business people, Big Data is coming for you.

Software that captures lots of data and uses it to make predictions has mostly been the province of engineers skilled in arcane databases and statisticians capable of developing complex algorithms. As the business gets bigger, however, software makers are domesticating their products in the hope they will prove attractive to a broader population.

[Cloudera](#), which offers a popular version of the open source database called Hadoop, released software on Wednesday that makes it possible to run queries from a more mainstream SQL programming language interface. SQL, thanks to its adoption by Oracle, Microsoft and others, is known to millions of business analysts.

“This enables us to talk to a whole other class of customer,” said Mike Olson, the chief executive of Cloudera. “The knock against Hadoop was that it is too complex.”

There is a reason for that. Hadoop is one of several so-called unstructured databases that were created at Yahoo and Google, after those two companies found they had previously unimaginable amounts of data about

PREVIOUS POST [Google Shifts Pitch for Its New Chromebooks](#)

NEXT POST [In Contest for Rescue Robots, Darpa Offers \\$2 Million Prize](#)

AROUND THE WEB »

THE NEXT WEB **Google says Maps redirect on Windows Phone was a product decision, and will be removed**

SCUTTLEBOT *News from the Web, annotated by our staff*

Google's Schmidt arrive in North Korea
REUTERS | From Mountain View to...errr, Pyongyang? - *Somini Sengupta*

AP provides sponsored tweets during electronics show
AP.ORG | The Associated Press is renting out its Twitter feed, with 1.5 million followers, to advertisers during C.E.S. - *Joshua Brustein*

A history of grieving
EDGE-ONLINE.COM | Meet the cult of gamers who want to ruin your day – just for kicks. - *Jenna Wortham*

A Million First Dates
THE ATLANTIC | Is online romance threatening monogamy? -

BLOOMBERG **HTC Posts Lowest Net Income in Eight Years After Revenue Drops**

And the XIXth Century...

Data changed, because of the digital/numeric revolution, see Gartner's 3V (Volume, Variety, Velocity), see [Gartner](#).

The New York Times
Sunday Review | The Opinion Pages

Search All NYTimes.com 

WORLD U.S. N.Y./REGION BUSINESS TECHNOLOGY SCIENCE HEALTH SPORTS OPINION ARTS STYLE TRAVEL JOBS REAL ESTATE AUTOS

NEWS ANALYSIS
The Age of Big Data
By STEVE LOHR
Published: February 11, 2012 | 82 Comments

GOOD with numbers? Fascinated by data? The sound you hear is opportunity knocking.

[Enlarge This Image](#)



Chad Hagen

Mo Zhou was snapped up by I.B.M. last summer, as a freshly minted Yale M.B.A., to join the technology company's fast-growing ranks of data consultants. They help businesses make sense of an explosion of data — Web traffic and social network comments, as well as software and sensors that monitor shipments, suppliers and customers — to guide decisions, trim costs and lift sales. "I've always had a love of numbers," says Ms. Zhou, whose job as a data analyst suits her skills.

To exploit the data flood, America will need many more like her. A report last year by the [McKinsey Global Institute](#), the research arm of the consulting firm, projected that the United States needs 140,000 to 190,000 more workers with "deep analytical" expertise and 1.5 million more data-literate managers, whether retrained or hired.

The impact of data abundance extends well beyond business. Justin Grimmer, for example, is one of the new breed of political scientists. A 28-year-old assistant

RECOMMEND
TWITTER
LINKEDIN
COMMENTS (82)
E-MAIL
PRINT
REPRINTS
SHARE

Log in to see what your friends are sharing on [Log In With Facebook](#)
nytimes.com. [Privacy Policy](#) | [What's This?](#)

What's Popular Now 

Felony Counts for 2 in Suicide of Bullied 12-Year-Old 

Senate Women Lead in Effort to Find Accord 

12 YEARS A SLAVE
[WATCH THE TRAILER](#)

Multimedia

Siera = $6.145 - 16.986 \times (S - 1.858 \times ((GB - FB - PU) \div PA) \times (((GB - FB - PU) \div PA)^2) + 1 \div PA) - 5.195 \times (BB \div PA) \times (($
whGraphic +/- term is a neg
Play (Data-Driven) Ball!

And the XIXth Century...

Business Intelligence, transversal approach

HOME PAGE TODAY'S PAPER VIDEO MOST POPULAR U.S. Edition ▼ Log In Register Now Help

The New York Times **Business Day** Search All NYTimes.com

WORLD U.S. N.Y. / REGION BUSINESS TECHNOLOGY SCIENCE HEALTH SPORTS OPINION ARTS STYLE TRAVEL JOBS REAL ESTATE AUTOS

Search Global DealBook Markets Economy Energy Media Personal Tech Small Business Your Money

UNBOXED

How Big Data Became So Big

By STEVE LOHR
Published: August 11, 2012

THIS has been the crossover year for Big Data — as a concept, as a term and, yes, as a marketing tool. Big Data has sprung from the confines of technology circles into the mainstream.

[Enlarge This Image](#)



Lloyd Miller for The New York Times

First, here are a few, well, data points: Big Data was a featured topic this year at the World Economic Forum in Davos, Switzerland, with a report titled "[Big Data, Big Impact.](#)" In March, the federal government announced \$200 million in research programs for Big Data computing.

Rick Smolan, creator of the "Day in the Life" photography series, has a new project in the works, called "The Human Face of Big Data." The New York Times has adopted the term in headlines like "[The Age of Big Data](#)" and "[Big Data on Campus.](#)" And a sure sign that Big Data has arrived

FACEBOOK
TWITTER
GOOGLE+
E-MAIL
SHARE
PRINT
REPRINTS

Log in to see what your friends are sharing on nytimes.com. [Log In With Facebook](#)
Privacy Policy | What's This?

What's Popular Now

Despite New Health Law, Some See Sharp Rise in Premiums  The Big Fail 

MOST E-MAILED **RECOMMENDED FOR YOU**

- OFF THE DRIBBLE
Stoudemire Commemorates Brother's Death
- CRITIC'S NOTEBOOK
The Rainbow That Follows 'Jersey Shore'
- TAKING NOTE
Opinion Report: Tax Reform
- THE LEARNING NETWORK
Fill-In | Trendy Spot Urges Tourists to Ride In and Spend, 'Gangnam Style'
- Major Companies Push the Limits of a Tax

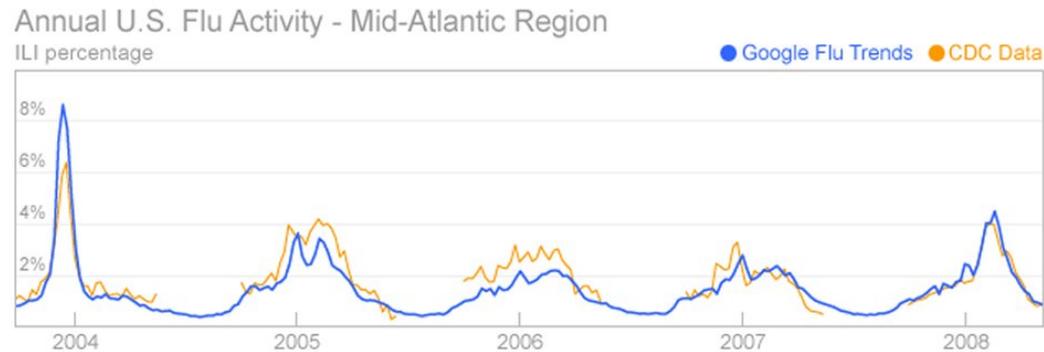
Add to Portfolio

[+ International Business Machines Corporation](#)

[Go to your Portfolio »](#)

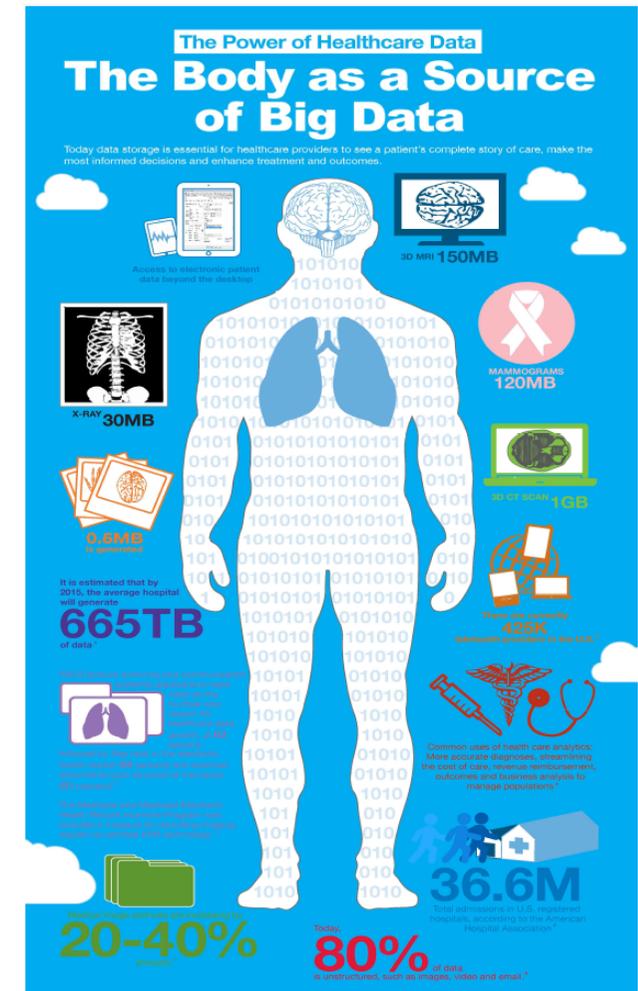
Big Data & (Health) Insurance

Example: popular application, Google Flu Trend



See also Lazer *et al.* (2014)

But much more can be done on an individual level.



Big Data & Computational Issues

parallel computing is a necessity*

CPU **Central Processing Unit**, the heart of the computer

RAM **Random Access Memory** non-persistent memory

HD **Hard Drive** persistent memory

Practical issues: CPU can be fast, but finite speed;

RAM is non persistent, fast but slow vs. HD is persistent, slow but big

How could we measure speed: Latency and performance

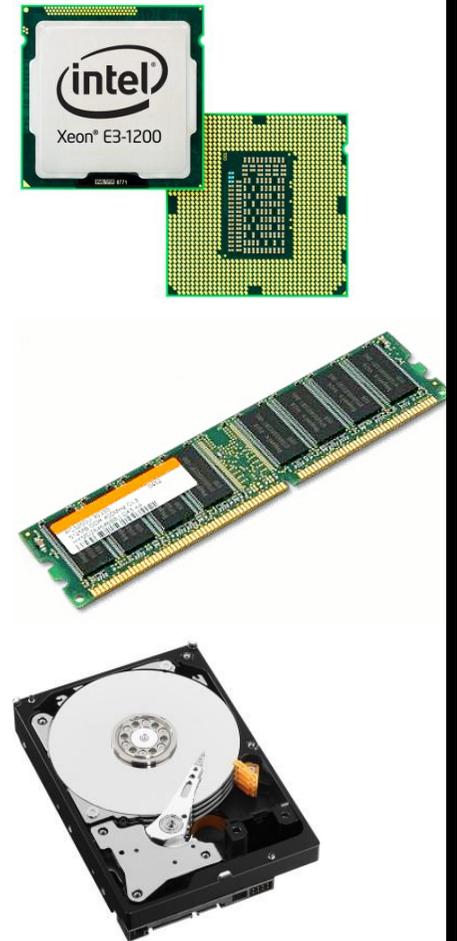
Latency is a time interval between the stimulation and response (e.g. 10ms to read the first bit)

Performance is the number of operations per second (e.g. 100Mb/sec)

Example Read one file of 100Mb \sim 1.01sec.

Example Read 150 files of 1b \sim 0.9sec.

* thanks to David Sibai for this section.



Big Data & Computational Issues

Standard PC :

CPU : 4 core, 1ns latency

RAM : 32 or 64 Gb, 100ns latency, 20Gb/sec

HD : 1 Tb, 10ms latency, 100Mo/sec

How long does it take ?

e.g. count *spaces* in a 2Tb text file

about $2 \cdot 10^{12}$ operations (comparaison)

File on the HD, 100Mb/sec $\sim 2 \cdot 10^4$ sec ~ 6 hours

COMPUTER SPECIALIST

DELL T3400

INTEL

CORE 2 DUO

320GB HDD

4GB RAM

256MB

GRAPHIC CARD



Big Data & Computational Issues

Why not parallelize ? between machines

Spread data on 10 blocks of 200Gb, each machine count spaces, then sum the 10 totals... should be 10 times faster.

Many machines connected, in a datacenter

Alternative: use more cores in the CPU (2, 4, 16 cores, e.g.)

A CPU is multi-tasks, and it could be possible to vectorize. E.g. summing n numbers takes $O(n)$ operations,

Example $a_1 + b_1, a_2 + b_2, \dots, a_n + b_n$ takes n nsec.

But it is possible to use SIMD (**single instruction multiple data**)

Example $\mathbf{a} + \mathbf{b} = (a_1, \dots, a_n) + (b_1, \dots, b_n)$ take 1 nsec.

Big Data & Computational Issues

Alternatives to standard PC material

Games from the 90s, more and more 3d viz, based on more and more computations

GPU **Graphical Processing Unit** that became GPGPU **General Purpose GPU**



Hundreds of small processors, slow, high specialized (and dedicated to simple computations)

Difficult to use (needs of computational skills) but more and more libraries

Complex and slow communication CPU - RAM - GPU

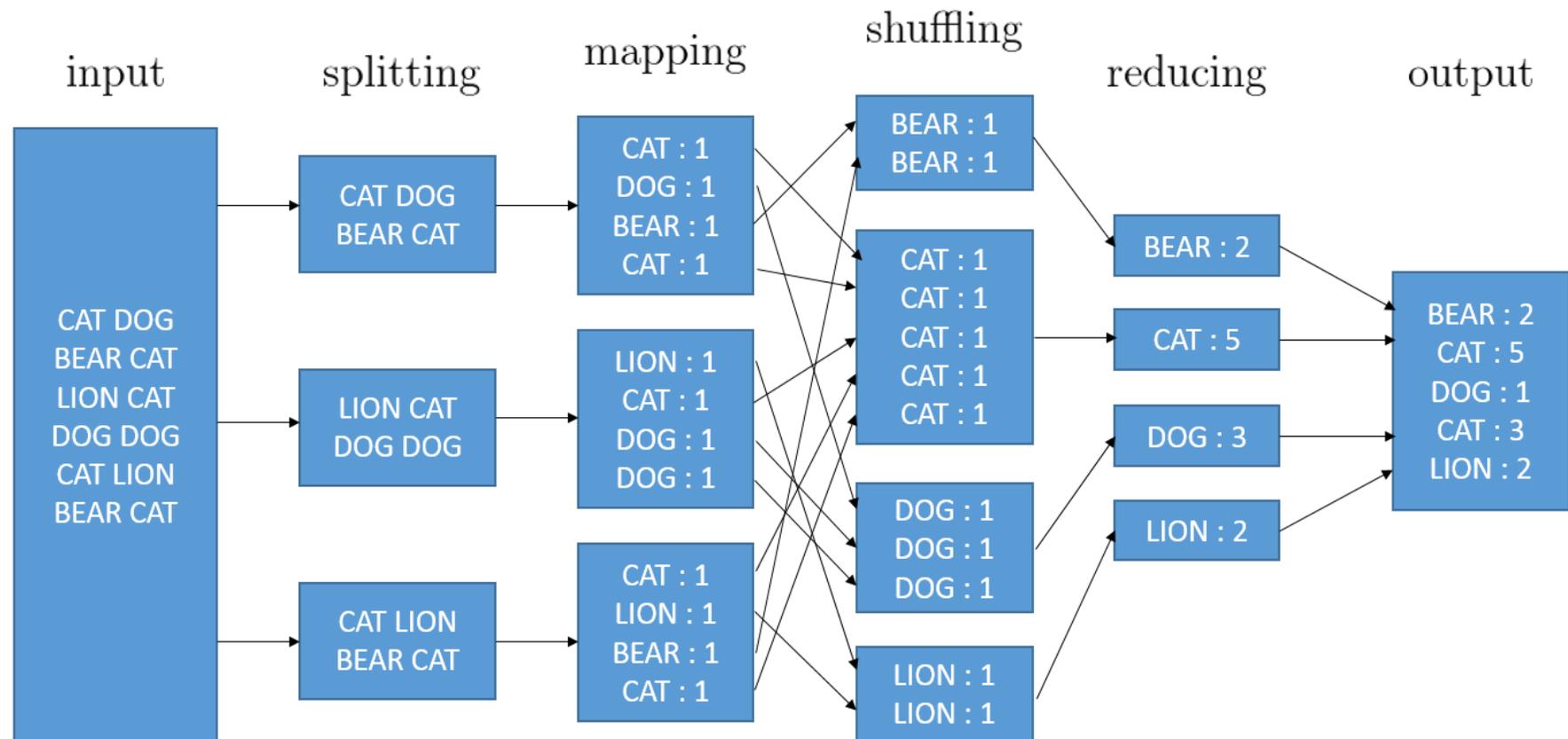
Sequential code is extremely slow, but highly parallelized

Interesting for Monte Carlo computations

E.g. pricing of **Variable Annuities**

Big Data & Computational Issues

A parallel algorithm is a computational strategy which divide a target computation into independent part, and assemble them so as to obtain the target computation. E.g. Counting words with MapReduce



Data, (deep) Learning & AI

What can we do with those data?

Part 2. Big Data and Statistical/Machine Learning



Statistical Learning and Philosophical Issues

From [Machine Learning and Econometrics](#), by Hal Varian :

“**Machine learning** use data to predict some variable as a function of other covariables,

- may, or may not, care about insight, importance, patterns
- may, or may not, care about inference (how y changes as some x change)

Econometrics use statistical methodes for prediction, inference and causal modeling of economic relationships

- hope for some sort of insight (inference is a goal)
- in particular, causal inference is goal for decision making.”

→ machine learning, ‘new tricks for econometrics’

Statistical Learning and Philosophical Issues

Remark machine learning can also learn from econometrics, especially with non i.i.d. data (time series and panel data)

Remark machine learning can help to get better predictive models, given good datasets. No use on several data science issues (e.g. selection bias).

non-supervised vs. supervised techniques

Non-Supervised and Supervised Techniques

Just \mathbf{x}_i 's, here, no y_i : unsupervised.

Use **principal components** to reduce dimension: we want d vectors $\mathbf{z}_1, \dots, \mathbf{z}_d$ such that

$$\mathbf{x}_i \sim \sum_{j=1}^d \omega_{i,j} \mathbf{z}_j \text{ or } \mathbf{X} \sim \mathbf{Z}\mathbf{\Omega}^T$$

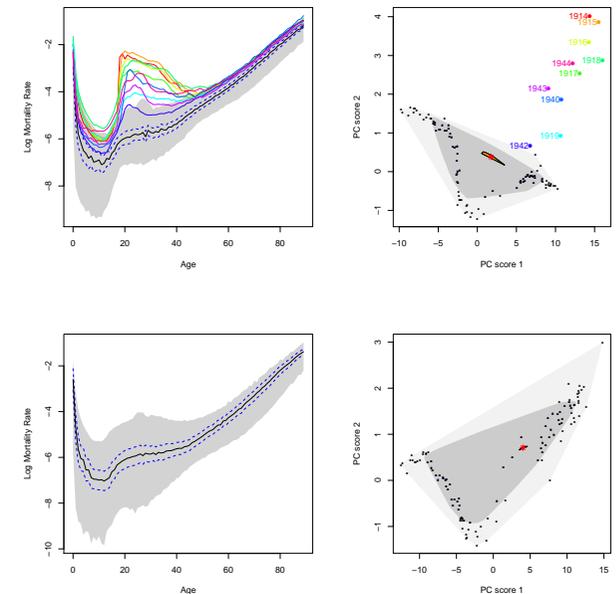
where $\mathbf{\Omega}$ is a $k \times d$ matrix, with $d < k$.

First Component is $\mathbf{z}_1 = \mathbf{X}\boldsymbol{\omega}_1$ where

$$\boldsymbol{\omega}_1 = \operatorname{argmax}_{\|\boldsymbol{\omega}\|=1} \left\{ \|\mathbf{X} \cdot \boldsymbol{\omega}\|^2 \right\} = \operatorname{argmax}_{\|\boldsymbol{\omega}\|=1} \left\{ \boldsymbol{\omega}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\omega} \right\}$$

Second Component is $\mathbf{z}_2 = \mathbf{X}\boldsymbol{\omega}_2$ where

$$\boldsymbol{\omega}_2 = \operatorname{argmax}_{\|\boldsymbol{\omega}\|=1} \left\{ \|\widetilde{\mathbf{X}}^{(1)} \cdot \boldsymbol{\omega}\|^2 \right\} \text{ where } \widetilde{\mathbf{X}}^{(1)} = \mathbf{X} - \underbrace{\mathbf{X}\boldsymbol{\omega}_1}_{\mathbf{z}_1} \boldsymbol{\omega}_1^T$$



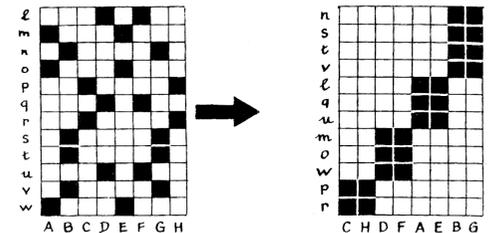
Unsupervised Techniques: Cluster Analysis

Data : $\{\mathbf{x}_i = (x_{1,i}, x_{2,i}), i = 1, \dots, n\}$

Distance matrix $D_{i,j} = D(\mathbf{x}_{c_i}, \mathbf{x}_{c_j})$

the distance is between clusters, not (only) individuals,

$$D(\mathbf{x}_{c_1}, \mathbf{x}_{c_2}) = \begin{cases} \min_{i \in c_1, j \in c_2} \{d(\mathbf{x}_i, \mathbf{x}_j)\} \\ d(\bar{\mathbf{x}}_{c_1}, \bar{\mathbf{x}}_{c_2}) \\ \max_{i \in c_1, j \in c_2} \{d(\mathbf{x}_i, \mathbf{x}_j)\} \end{cases}$$



for some (standard) distance d , e.g. Euclidean (ℓ_2), Manhattan (ℓ_1), Jaccard, etc.
See also [Bertin \(1967\)](#).

Unsupervised Techniques

Data : $\{\mathbf{x}_i = (x_{1,i}, x_{2,i}), i = 1, \dots, n\}$

\mathbf{x}_i 's are observations from i.i.d random variables \mathbf{X}_i
with distribution $F_{\mathbf{p}, \boldsymbol{\theta}}$,

$$F_{\mathbf{p}, \boldsymbol{\theta}}(\mathbf{x}) = \underbrace{p_1 \cdot F_{\boldsymbol{\theta}_1}(\mathbf{x})}_{\text{Cluster 1}} + \underbrace{p_2 \cdot F_{\boldsymbol{\theta}_2}(\mathbf{x})}_{\text{Cluster 2}} + \dots$$

E.g. $F_{\boldsymbol{\theta}_k}$ is the c.d.f. of a $\mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ distribution.

Unsupervised Techniques

Data : $\{\mathbf{x}_i = (x_{1,i}, x_{2,i}), i = 1, \dots, n\}$

iterative procedure:

1. start with k points $\mathbf{z}_1, \dots, \mathbf{z}_k$
2. cluster c_j are $\{d(\mathbf{x}_i, \mathbf{z}_j) \leq d(\mathbf{x}_i, \mathbf{z}_{j'}), j' \neq j\}$
3. $\mathbf{z}_j = \bar{\mathbf{x}}_{c_j}$

See Steinhaus (1957) or Lloyd (1957)

But **curse of dimensionality**, unhelpful in high dimension

Datamining, Explantory Analysis, Regression, Statistical Learning, Predictive Modeling, etc

In statistical learning, data are approached with little priori information.

In regression analysis, see [Cook & Weisberg \(1999\)](#)

The primary goal in a regression analysis is to understand, as far as possible with the available data, how the conditional distribution of the response y varies across subpopulations determined by the possible values of the predictor or predictors. Since this is the central idea, it will be helpful to have a conve-

i.e. we would like to get the distribution of the response variable Y conditioning on one (or more) predictors \mathbf{X} .

Consider a regression model, $y_i = m(\mathbf{x}_i) + \varepsilon_i$, where ε_i 's are i.i.d. $\mathcal{N}(0, \sigma^2)$, possibly linear $y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i$, where ε_i 's are (somehow) unpredictable.

Machine Learning and ‘Statistics’

Machine learning and statistics seem to be very similar, they share the same goals—they both focus on data modeling—but their methods are affected by their cultural differences.

“The goal for a statistician is to predict an interaction between variables with some degree of certainty (we are never 100% certain about anything). Machine learners, on the other hand, want to build algorithms that predict, classify, and cluster with the most accuracy, see [Why a Mathematician, Statistician & Machine Learner Solve the Same Problem Differently](#)

Machine learning methods are about algorithms, more than about asymptotic statistical properties.

Validation is not based on mathematical properties, but on properties out of sample: we must use a [training sample](#) to train (estimate) model, and a [testing sample](#) to compare algorithms (hold out technique).

Goldilock Principle: the Mean-Variance Tradeoff

In statistics and in machine learning, there will be **parameters** and **meta-parameters** (or **tunning parameters**). The first ones are estimated, the second ones should be chosen.

See **Hill estimator** in **extreme value** theory. X has a Pareto distribution - with index ξ - above some threshold u if

$$\mathbb{P}[X > x | X > u] = \left(\frac{u}{x}\right)^{\frac{1}{\xi}} \text{ for } x > u.$$

Given a sample \mathbf{x} , consider the Pareto-QQ plot, i.e. the scatterplot

$$\left\{ -\log \left(1 - \frac{i}{n+1} \right), \log x_{i:n} \right\}_{i=n-k, \dots, n}$$

for points exceeding $X_{n-k:n}$. The slope is ξ , i.e.

$$\log X_{n-i+1:n} \approx \log X_{n-k:n} + \xi \left(-\log \frac{i}{n+1} - \log \frac{n+1}{k+1} \right)$$

Goldilock Principle: the Mean-Variance Tradeoff

Hence, consider estimator

$$\hat{\xi}_k = \frac{1}{k} \sum_{i=0}^{k-1} \log x_{n-i:n} - \log x_{n-k:n}.$$

k is the number of large observations, in the upper tail.

Standard mean-variance tradeoff,

- k large: bias too large, variance too small
- k small: variance too large, bias too small

Goldilock Principle: the Mean-Variance Tradeoff

Same holds in **kernel regression**, with bandwidth h (length of neighborhood)

$$\hat{m}_h(x) = \frac{\sum_{i=1}^n K_h(x - x_i) y_i}{\sum_{i=1}^n K_h(x - x_i)}$$

since

$$\mathbb{E}(Y|X = x) = \int \frac{f(x, y) \cdot y}{f(x)} dy$$

Standard **mean-variance tradeoff**,

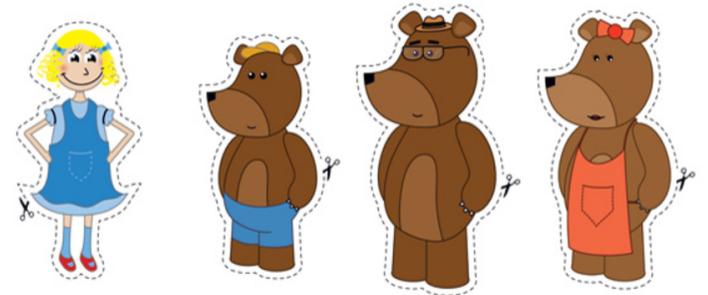
- h large: bias too large, variance too small
- h small: variance too large, bias too small

Goldilock Principle: the Mean-Variance Tradeoff

More generally, we estimate $\hat{\theta}_h$ or $\hat{m}_h(\cdot)$

Use the **mean squared error** for $\hat{\theta}_h$

$$\mathbb{E} \left[\left(\theta - \hat{\theta}_h \right)^2 \right]$$



or **mean integrated squared error** $\hat{m}_h(\cdot)$,

$$\mathbb{E} \left[\int (m(\mathbf{x}) - \hat{m}_h(\mathbf{x}))^2 d\mathbf{x} \right]$$

In statistics, derive an asymptotic expression for these quantities, and find h^* that minimizes those.

Goldilock Principle: the Mean-Variance Tradeoff

For kernel regression, the MISE can be approximated by

$$\frac{h^4}{4} \left(\int \mathbf{x}^\top \mathbf{x} K(\mathbf{x}) d\mathbf{x} \right)^2 \int \left(m''(\mathbf{x}) + 2m'(\mathbf{x}) \frac{f'(\mathbf{x})}{f(\mathbf{x})} \right) d\mathbf{x} + \frac{1}{nh} \sigma^2 \int K^2(\mathbf{x}) d\mathbf{x} \int \frac{d\mathbf{x}}{f(\mathbf{x})}$$

where f is the density of \mathbf{x} 's. Thus the optimal h is

$$h^* = n^{-\frac{1}{5}} \left(\frac{\sigma^2 \int K^2(\mathbf{x}) d\mathbf{x} \int \frac{d\mathbf{x}}{f(\mathbf{x})}}{\left(\int \mathbf{x}^\top \mathbf{x} K(\mathbf{x}) d\mathbf{x} \right)^2 \int \left(\int m''(\mathbf{x}) + 2m'(\mathbf{x}) \frac{f'(\mathbf{x})}{f(\mathbf{x})} \right)^2 d\mathbf{x}} \right)^{\frac{1}{5}}$$

(hard to get a simple rule of thumb... up to a constant, $h^* \sim n^{-\frac{1}{5}}$)

Use bootstrap, or cross-validation to get an optimal h

Randomization is too important to be left to chance!

Bootstrap (resampling) algorithm is very important (nonparametric monte carlo)

→ data (and not model) driven algorithm

Randomization is too important to be left to chance!

Consider some sample $\mathbf{x} = (x_1, \dots, x_n)$ and some statistics $\hat{\theta}$. Set $\hat{\theta}_n = \hat{\theta}(\mathbf{x})$

Jackknife used to reduce bias: set $\hat{\theta}_{(-i)} = \hat{\theta}(\mathbf{x}_{(-i)})$, and $\tilde{\theta} = \frac{1}{n} \sum_{i=1}^n \hat{\theta}_{(-i)}$

If $\mathbb{E}(\hat{\theta}_n) = \theta + O(n^{-1})$ then $\mathbb{E}(\tilde{\theta}_n) = \theta + O(n^{-2})$.

See also **leave-one-out** cross validation, for $\hat{m}(\cdot)$

$$\text{mse} = \frac{1}{n} \sum_{i=1}^n [y_i - \hat{m}_{(-i)}(x_i)]^2$$

Bootstrap estimate is based on bootstrap samples: set $\hat{\theta}_{(b)} = \hat{\theta}(\mathbf{x}_{(b)})$, and

$\tilde{\theta} = \frac{1}{n} \sum_{i=1}^n \hat{\theta}_{(b)}$, where $\mathbf{x}_{(b)}$ is a vector of size n , where values are drawn from

$\{x_1, \dots, x_n\}$, with replacement. And then use the law of large numbers...

See **Efron (1979)**.

Hold-Out, Cross Validation, Bootstrap

Hold-out: Split $\{1, \dots, n\}$ into T (training) and V (validation)

Train the model on $\{(y_i, \mathbf{x}_i), i \in T\}$ and compute

$$\hat{R} = \frac{1}{\#(V)} \sum_{i \in V} \ell(y_i, \hat{m}(\mathbf{x}_i))$$

k -fold cross validation: Split $\{1, \dots, n\}$ into I_1, \dots, I_k . Set $I_{\bar{j}} = \{1, \dots, n\} \setminus I_j$

Train model on $I_{\bar{j}}$ and compute

$$\hat{R} = \frac{1}{k} \sum_j R_j \text{ where } R_j = \frac{k}{n} \sum_{i \in I_j} \ell(y_i, \hat{m}_{\bar{j}}(\mathbf{x}_i))$$

Hold-Out, Cross Validation, Bootstrap

Leave-one-out bootstrap: generate I_1, \dots, I_B bootstrapped samples from $\{1, \dots, n\}$

set $n_i = \mathbf{1}_{i \notin I_1} + \dots + \mathbf{1}_{i \notin I_B}$

$$\hat{R} = \frac{1}{n} \sum_{i=1}^n \frac{1}{n_i} \sum_{b: i \notin I_b} \ell(y_i, \hat{m}_b(\mathbf{x}_i))$$

Remark Probability that i th row is not selection $(1 - n^{-1})^n \rightarrow e^{-1} \sim 36.8\%$,
cf training / validation samples (2/3-1/3)

Statistical Learning and Philosophical Issues

From (y_i, \mathbf{x}_i) , there are different stories behind, see [Freedman \(2005\)](#)

- the **causal story** : $x_{j,i}$ is usually considered as independent of the other covariates $x_{k,i}$. For all possible \mathbf{x} , that value is mapped to $m(\mathbf{x})$ and a noise is attached, ε . The goal is to recover $m(\cdot)$, and the residuals are just the difference between the response value and $m(\mathbf{x})$.
- the **conditional distribution story** : for a linear model, we usually say that Y given $\mathbf{X} = \mathbf{x}$ is a $\mathcal{N}(m(\mathbf{x}), \sigma^2)$ distribution. $m(\mathbf{x})$ is then the conditional mean. Here $m(\cdot)$ is assumed to really exist, but no causal assumption is made, only a conditional one.
- the **explanatory data story** : there is no model, just data. We simply want to summarize information contained in \mathbf{x} 's to get an accurate summary, close to the response (i.e. $\min\{\ell(\mathbf{y}_i, m(\mathbf{x}_i))\}$) for some loss function ℓ .

Machine Learning vs. Statistical Modeling

In **machine learning**, given some dataset (\mathbf{x}_i, y_i) , solve

$$\hat{m}(\cdot) = \operatorname{argmin}_{m(\cdot) \in \mathcal{F}} \left\{ \sum_{i=1}^n \ell(y_i, m(\mathbf{x}_i)) \right\}$$

for some **loss functions** $\ell(\cdot, \cdot)$.

In **statistical modeling**, given some probability space $(\Omega, \mathcal{A}, \mathbb{P})$, assume that y_i are realization of i.i.d. variables Y_i (given $\mathbf{X}_i = \mathbf{x}_i$) with distribution F_i . Then solve

$$\hat{m}(\cdot) = \operatorname{argmax}_{m(\cdot) \in \mathcal{F}} \{ \log \mathcal{L}(m(\mathbf{x}); \mathbf{y}) \} = \operatorname{argmax}_{m(\cdot) \in \mathcal{F}} \left\{ \sum_{i=1}^n \log f(y_i; m(\mathbf{x}_i)) \right\}$$

where $\log \mathcal{L}$ denotes the **log-likelihood**.

Computational Aspects: Optimization

Econometrics, Statistics and Machine Learning rely on the same object:
optimization routines.

A gradient descent/ascent algorithm

A stochastic algorithm

Loss Functions

Fitting criteria are based on **loss functions** (also called **cost functions**). For a quantitative response, a popular one is the quadratic loss,

$$\ell(y, m(\mathbf{x})) = [y - m(\mathbf{x})]^2.$$

Recall that

$$\left\{ \begin{array}{l} \mathbb{E}(Y) = \operatorname{argmin}_{m \in \mathbb{R}} \{ \|Y - m\|_{\ell_2}^2 \} = \operatorname{argmin}_{m \in \mathbb{R}} \{ \mathbb{E}([Y - m]^2) \} \\ \operatorname{Var}(Y) = \min_{m \in \mathbb{R}} \{ \mathbb{E}([Y - m]^2) \} = \mathbb{E}([Y - \mathbb{E}(Y)]^2) \end{array} \right.$$

The empirical version is

$$\left\{ \begin{array}{l} \bar{y} = \operatorname{argmin}_{m \in \mathbb{R}} \left\{ \sum_{i=1}^n \frac{1}{n} [y_i - m]^2 \right\} \\ s^2 = \min_{m \in \mathbb{R}} \left\{ \sum_{i=1}^n \frac{1}{n} [y_i - m]^2 \right\} = \sum_{i=1}^n \frac{1}{n} [y_i - \bar{y}]^2 \end{array} \right.$$

Loss Functions

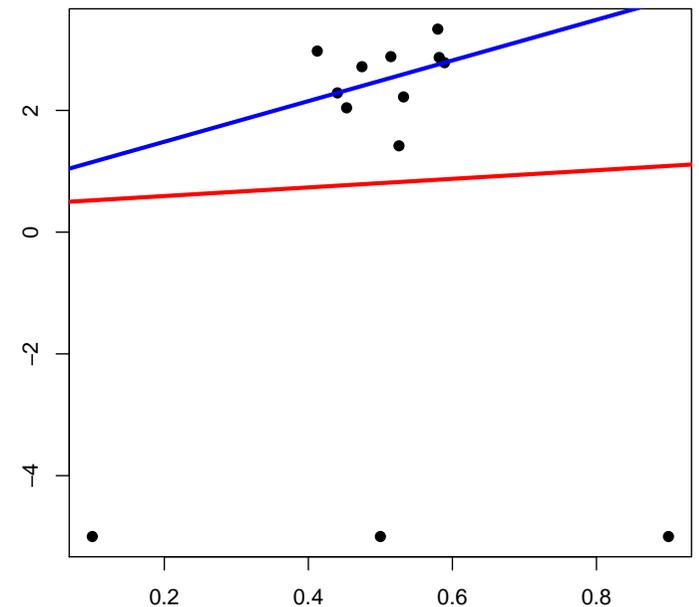
Remark $\text{median}(\mathbf{y}) = \underset{m \in \mathbb{R}}{\text{argmin}} \left\{ \sum_{i=1}^n \frac{1}{n} |y_i - m| \right\}$

Quadratic loss function $\ell(a, b)^2 = (a - b)^2$,

$$\sum_{i=1}^n (y_i - \mathbf{x}^T \boldsymbol{\beta})^2 = \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|_{\ell_2}^2$$

Absolute loss function $\ell(a, b) = |a - b|$

$$\sum_{i=1}^n |y_i - \mathbf{x}^T \boldsymbol{\beta}| = \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|_{\ell_1}$$



Loss Functions

Quadratic loss function $\ell_2(x, y)^2 = (x - y)^2$,

Absolute loss function $\ell_1(x, y) = |x - y|$

Quantile loss function $\ell_\tau(x, y) = |(x - y)(\tau - \mathbf{1}_{x \leq y})|$

Huber loss function

$$\ell_\tau(x, y) = \begin{cases} \frac{1}{2}(x - y)^2 & \text{for } |x - y| \leq \tau, \\ \tau |x - y| - \frac{1}{2}\tau^2 & \text{otherwise.} \end{cases}$$

i.e. quadratic when $|x - y| \leq \tau$ and linear otherwise.

Loss Functions

For classification: **misclassification** loss function

$$\ell(x, y) = \mathbf{1}_{x \neq y} \text{ or } \ell(x, y) = \mathbf{1}_{\text{sign}(x) \neq \text{sign}(y)}$$

$$\ell_{\tau}(x, y) = \tau \mathbf{1}_{\text{sign}(x) < 0, \text{sign}(y) > 0} + [1 - \tau] \mathbf{1}_{\text{sign}(x) > 0, \text{sign}(y) < 0}$$

For $\{-1, +1\}$ classes,

Hinge loss ('maximum-margin' classification) $\ell(x, y) = (1 - xy)_+$

Logistic/log loss $\ell(x, y) = \log[1 + e^{-xy}]$

Squared Loss $\ell(x, y) = [x - y]^2 = [1 - xy]^2$

Linear Predictors

In the linear model, least square estimator yields

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \underbrace{\mathbf{X}[\mathbf{X}^T\mathbf{X}]^{-1}\mathbf{X}^T}_{\mathbf{H}}\mathbf{Y}$$

We have a **linear predictor** if the fitted value \hat{y} at point \mathbf{x} can be written

$$\hat{y} = \hat{m}(\mathbf{x}) = \sum_{i=1}^n \mathbf{S}_{\mathbf{x},i} y_i = \mathbf{S}_{\mathbf{x}}^T \mathbf{y}$$

where $\mathbf{S}_{\mathbf{x}}$ is some vector of weights (called **smoother vector**), related to a $n \times n$ smoother matrix,

$$\hat{\mathbf{y}} = \mathbf{S}\mathbf{y}$$

where prediction is done at points \mathbf{x}_i 's.

Degrees of Freedom and Model Complexity

E.g.

$$\mathbf{S}_x = \mathbf{X}[\mathbf{X}^\top \mathbf{X}]^{-1} \mathbf{x}$$

that is related to the hat matrix, $\hat{\mathbf{y}} = \mathbf{H}\mathbf{y}$.

Note that

$$T = \frac{\|\mathbf{S}\mathbf{Y} - \mathbf{H}\mathbf{Y}\|}{\text{trace}([\mathbf{S} - \mathbf{H}]^\top [\mathbf{S} - \mathbf{H}])}$$

can be used to test a linear assumption: if the model is linear, then T has a Fisher distribution.

In the context of linear predictors, $\text{trace}(\mathbf{S})$ is usually called **equivalent number of parameters** and is related to n - effective degrees of freedom (as in [Ruppert et al. \(2003\)](#)).

Model Evaluation

In linear models, the R^2 is defined as the proportion of the variance of the the response y that can be obtained using the predictors.

But maximizing the R^2 usually yields **overfit** (or **unjustified optimism** in Berk (2008)).

In linear models, consider the adjusted R^2 ,

$$\bar{R}^2 = 1 - [1 - R^2] \frac{n - 1}{n - p - 1}$$

where p is the number of parameters (or more generally $\text{trace}(\mathbf{S})$).

Model Evaluation

Alternatives are based on the Akaike Information Criterion (**AIC**) and the Bayesian Information Criterion (**BIC**), based on a penalty imposed on some criteria (the logarithm of the variance of the residuals),

$$AIC = \log \left(\frac{1}{n} \sum_{i=1}^n [y_i - \hat{y}_i]^2 \right) + \frac{2p}{n}$$

$$BIC = \log \left(\frac{1}{n} \sum_{i=1}^n [y_i - \hat{y}_i]^2 \right) + \frac{\log(n)p}{n}$$

In a more general context, replace p by $\text{trace}(\mathbf{S})$

Goodhart's Law

‘when a measure becomes a target, it ceases to be a good measure’, by Charles Goodhart



Occam's Razor

Between two models which explain as well the data, choose the simplest one



Machine Learning: usually need to tradeoff between the training error and model complexity

$$\hat{m} = \underset{m}{\operatorname{argmin}} \{ \ell(Y, m(\mathbf{X})) + \Omega(m) \}$$

Model Evaluation

One can also consider the expected prediction error (with a probabilistic model)

$$\mathbb{E}[\ell(Y, \hat{m}(\mathbf{X}))]$$

We cannot claim (using the law of large number) that

$$\frac{1}{n} \sum_{i=1}^n \ell(y_i, \hat{m}(\mathbf{x}_i)) \stackrel{a.s.}{\not\rightarrow} \mathbb{E}[\ell(Y, m(\mathbf{X}))]$$

since \hat{m} depends on (y_i, \mathbf{x}_i) 's.

Natural option : use two (random) samples, a **training** one and a **validation** one.

Alternative options, use cross-validation, leave-one-out or k -fold.

Underfit / Overfit and Variance - Mean Tradeoff

Goal in predictive modeling: reduce uncertainty in our predictions.

Need more data to get a better knowledge.

Unfortunately, reducing the error of the prediction on a dataset does not generally give a good **generalization** performance

→ need a training and a validation dataset

Overfit, Training vs. Validation and Complexity

complexity \longleftrightarrow polynomial degree

Overfit, Training vs. Validation and Complexity

complexity \longleftrightarrow number of neighbors (k)

Themes in Data Science

Predictive Capability we want here to have a model that predict well for new observations

Bias-Variance Tradeoff A very smooth prediction has less variance, but a large bias. We need to find a good balance between the bias and the variance

Loss Functions In machine learning, goodness of fit is discussed based on disparities between predicted values, and observed one, based on some loss function

Tuning or Meta Parameters Choice will be made in terms of tuning parameters

Interpretability Does it matter to have a good model if we cannot interpret it ?

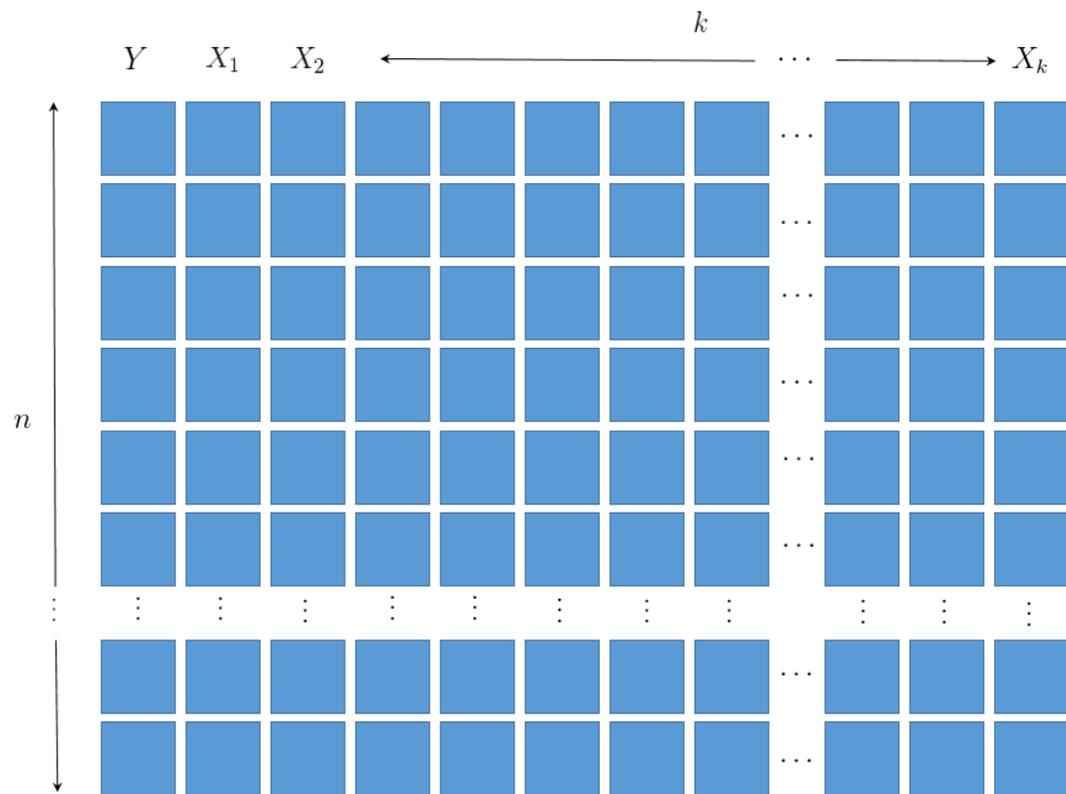
Coding Issues Most of the time, there are no analytical expression, just an algorithm that should converge to some (possibly) optimal value

Data Data collection is a crucial issue (but will not be discussed here)

Scalability Issues

Dealing with **big** (or **massive**) datasets, large number of observations (n) and/or large number of predictors (features or covariates, k).

Ability to parallelize algorithms might be important (map-reduce).



n can be large, but limited
(portfolio size)

large **variety** k

large **volume** nk

→ Feature Engineering

Part 3. Application to Classification



Econometric Based Models in Actuarial Science

Consider an i.i.d. sample $\{y_1, \dots, y_n\}$ with $y_i \in \{0, 1\}$,

$$\mathbb{P}(Y_i = y_i) = \pi^{y_i} [1 - \pi]^{1-y_i}, \text{ with } y_i \in \{0, 1\}.$$

where $\pi \in [0, 1]$, so that $\mathbb{P}(Y_i = 1) = \pi$ and $\mathbb{P}(Y_i = 0) = 1 - \pi$.

The likelihood is

$$\mathcal{L}(\pi; \mathbf{y}) = \prod_{i=1}^n \mathbb{P}(Y_i = y_i) = \prod_{i=1}^n \pi^{y_i} [1 - \pi]^{1-y_i}$$

and the **log-likelihood** is

$$\log \mathcal{L}(\pi; \mathbf{y}) = \sum_{i=1}^n y_i \log[\pi] + (1 - y_i) \log[1 - \pi]$$

The first order condition is

$$\frac{\partial \log \mathcal{L}(\pi; \mathbf{y})}{\partial \pi} = \sum_{i=1}^n \frac{y_i}{\pi} - \frac{1 - y_i}{1 - \pi} = 0, \text{ i.e. } \pi^* = \bar{y}.$$

Econometric Based Models in Actuarial Science

Assume that $\mathbb{P}(Y_i = 1) = \pi_i$,

$$\text{logit}(\pi_i) = \mathbf{X}'_i \boldsymbol{\beta}, \text{ where } \text{logit}(\pi_i) = \log \left(\frac{\pi_i}{1 - \pi_i} \right),$$

or

$$\pi_i = \text{logit}^{-1}(\mathbf{X}'_i \boldsymbol{\beta}) = \frac{\exp[\mathbf{X}'_i \boldsymbol{\beta}]}{1 + \exp[\mathbf{X}'_i \boldsymbol{\beta}]}.$$

The log-likelihood is

$$\log \mathcal{L}(\boldsymbol{\beta}) = \sum_{i=1}^n y_i \log(\pi_i) + (1 - y_i) \log(1 - \pi_i) = \sum_{i=1}^n y_i \log(\pi_i(\boldsymbol{\beta})) + (1 - y_i) \log(1 - \pi_i(\boldsymbol{\beta}))$$

and the first order conditions are **solved numerically**

$$\frac{\partial \log \mathcal{L}(\boldsymbol{\beta})}{\partial \beta_k} = \sum_{i=1}^n X_{k,i} [y_i - \pi_i(\boldsymbol{\beta})] = 0.$$

Logistic Classification

It is a linear classifier since we have a linear separation between the ●'s and the ●'s.

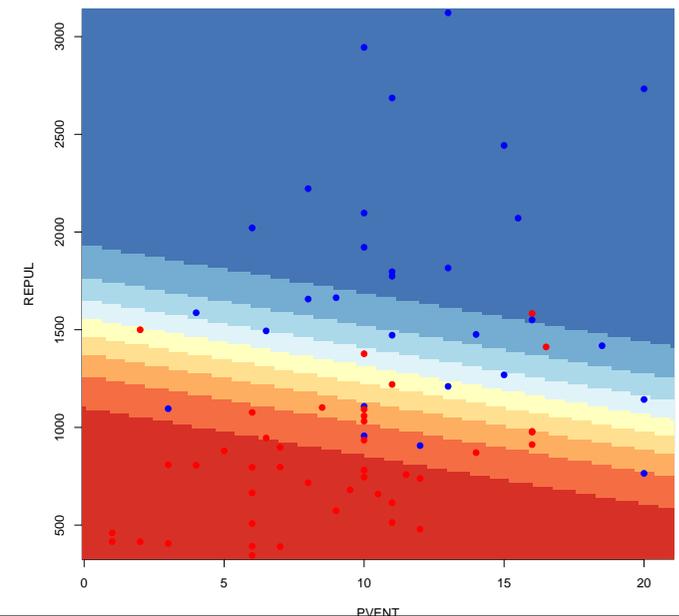
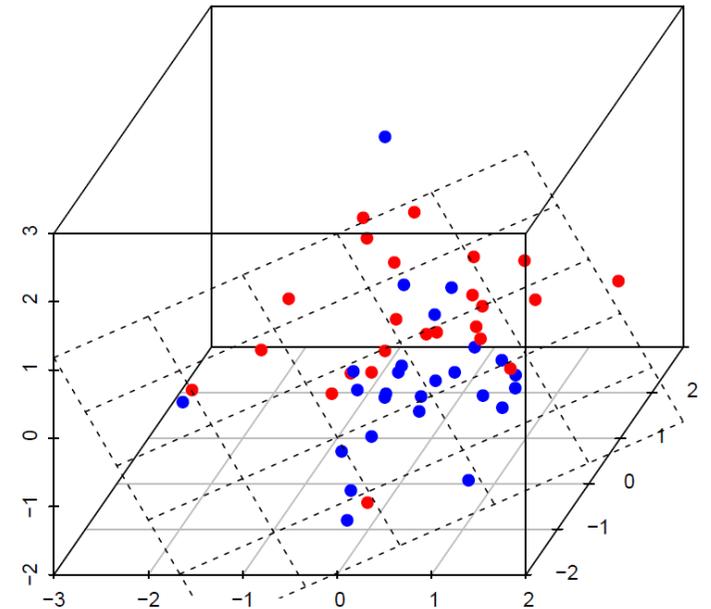
Let $m(\mathbf{x}) = \mathbb{E}(Y|\mathbf{X} = \mathbf{x})$.

With a logistic regression, we can get a prediction

$$\hat{m}(\mathbf{x}) = \frac{\exp[\mathbf{x}^\top \hat{\boldsymbol{\beta}}]}{1 + \exp[\mathbf{x}^\top \hat{\boldsymbol{\beta}}]}$$

Is that the 'best' model we can get from the data?

What if n and/or k are very large?



Why a Logistic and not a Probit Regression?

Bliss (1934) suggested a model such that

$$\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x}) = H(\mathbf{x}^\top \boldsymbol{\beta}) \text{ where } H(\cdot) = \Phi(\cdot)$$

the c.d.f. of the $\mathcal{N}(0, 1)$ distribution. This is the **probit** model.

This yields a latent model, $y_i = \mathbf{1}(y_i^* > 0)$ where

$$y_i^* = \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i \text{ is a nonobservable score.}$$

In the logistic regression, we model the **odds ratio**,

$$\frac{\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x})}{\mathbb{P}(Y \neq 1 | \mathbf{X} = \mathbf{x})} = \exp[\mathbf{x}^\top \boldsymbol{\beta}]$$

$$\mathbb{P}(Y = 1 | \mathbf{X} = \mathbf{x}) = H(\mathbf{x}^\top \boldsymbol{\beta}) \text{ where } H(\cdot) = \frac{\exp[\cdot]}{1 + \exp[\cdot]}$$

which is the c.d.f. of the **logistic** variable, see **Verhulst (1845)**

Table 3.2 Transformation of percentages to probits

%	0	1	2	3	4	5	6	7	8	9
0	—	2.67	2.95	3.12	3.25	3.36	3.45	3.52	3.59	3.66
10	3.72	3.77	3.82	3.87	3.92	3.96	4.01	4.05	4.08	4.12
20	4.16	4.19	4.23	4.26	4.29	4.33	4.36	4.39	4.42	4.45
30	4.48	4.50	4.53	4.56	4.59	4.61	4.64	4.67	4.69	4.72
40	4.75	4.77	4.80	4.82	4.85	4.87	4.90	4.92	4.95	4.97
50	5.00	5.03	5.05	5.08	5.10	5.13	5.15	5.18	5.20	5.23
60	5.25	5.28	5.31	5.33	5.36	5.39	5.41	5.44	5.47	5.50
70	5.52	5.55	5.58	5.61	5.64	5.67	5.71	5.74	5.77	5.81
80	5.84	5.88	5.92	5.95	5.99	6.04	6.08	6.13	6.18	6.23
90	6.28	6.34	6.41	6.48	6.55	6.64	6.75	6.88	7.05	7.33
—	0.0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9
99	7.33	7.37	7.41	7.46	7.51	7.58	7.65	7.75	7.88	8.09

Soit p la population : représentons par dp l'accroissement infiniment petit qu'elle reçoit pendant un temps infiniment court dt . Si la population croissait en progression géométrique, nous aurions l'équation $\frac{dp}{dt} = mp$. Mais comme la vitesse d'accroissement de la population est retardée par l'augmentation même du nombre des habitants, nous devons retrancher de mp une fonction inconnue de p ; de manière que la formule à intégrer deviendra

$$\frac{dp}{dt} = mp - \varphi(p).$$

L'hypothèse la plus simple que l'on puisse faire sur la forme de la fonction φ , est de supposer $\varphi(p) = np^2$. On trouve alors pour intégrale de l'équation ci-dessus

$$t = \frac{1}{m} [\log p - \log(m - np)] + \text{constante},$$

et il suffira de trois observations pour déterminer les deux coefficients constants m et n et la constante arbitraire.

En résolvant la dernière équation par rapport à p , il vient

$$p = \frac{mp' e^{mt}}{np' e^{mt} + m - np'} \dots \dots (1)$$

en désignant par p' la population qui répond à $t = 0$, et par e la base des logarithmes népériens. Si l'on fait $t = \infty$, on voit que la valeur de p correspondante est $P = \frac{m}{n}$. Telle est donc la limite supérieure de la population.

Predictive Classifier

To go from a score to a class:

if $s(\mathbf{x}) > s$, then $\hat{Y}(\mathbf{x}) = 1$ and $s(\mathbf{x}) \leq s$, then $\hat{Y}(\mathbf{x}) = 0$

Plot $TP(s) = \mathbb{P}[\hat{Y} = 1|Y = 1]$ against $FP(s) = \mathbb{P}[\hat{Y} = 1|Y = 0]$

Comparing Classifiers: Accuracy and Kappa

Kappa statistic κ compares an Observed Accuracy with an Expected Accuracy (random chance), see [Landis & Koch \(1977\)](#).

	$Y = 0$	$Y = 1$	
$\hat{Y} = 0$	TN	FN	TN+FN
$\hat{Y} = 1$	FP	TP	FP+TP
	TN+FP	FN+TP	n

See also Observed and Random Confusion Tables

	$Y = 0$	$Y = 1$	
$\hat{Y} = 0$	25	3	28
$\hat{Y} = 1$	4	39	43
	29	42	71

	$Y = 0$	$Y = 1$	
$\hat{Y} = 0$	11.44	16.56	28
$\hat{Y} = 1$	17.56	25.44	43
	29	42	71

$$\text{total accuracy} = \frac{TP + TN}{n} \sim 90.14\%$$

$$\text{random accuracy} = \frac{[TN + FP] \cdot [TP + FN] + [TP + FP] \cdot [TN + FN]}{n^2} \sim 51.93\%$$

$$\kappa = \frac{\text{total accuracy} - \text{random accuracy}}{1 - \text{random accuracy}} \sim 79.48\%$$

On Model Selection

Consider predictions obtained from a linear model and a nonlinear model, either on the training sample, or on a validation sample,

Penalization and Support Vector Machines

SVMs were developed in the 90's based on previous work, from [Vapnik & Lerner \(1963\)](#), see also [Vailant \(1984\)](#).

Assume that points are [linearly separable](#), i.e. there is ω and b such that

$$Y = \begin{cases} +1 & \text{if } \omega^\top \mathbf{x} + b > 0 \\ -1 & \text{if } \omega^\top \mathbf{x} + b < 0 \end{cases}$$

Problem: infinite number of solutions, need a [good](#) one, that separate the data, (somehow) far from the data.

maximize the distance s.t. $H_{\omega,b}$ separates ± 1 points, i.e.

$$\min \left\{ \frac{1}{2} \omega^\top \omega \right\} \text{ s.t. } Y_i(\omega^\top \mathbf{x}_i + b) \geq 1, \forall i.$$

Penalization and Support Vector Machines

Define **support vectors** as observations such that

$$|\omega^T \mathbf{x}_i + b| = 1$$

The margin is the distance between hyperplanes defined by support vectors. The distance from support vectors to $H_{\omega,b}$ is $\|\omega\|^{-1}$

Now, what about the **non-separable case**?

Here, we **cannot** have $y_i(\omega^T \mathbf{x}_i + b) \geq 1 \forall i$.

Penalization and Support Vector Machines

Thus, introduce *slack variables*,

$$\begin{cases} \omega^\top \mathbf{x}_i + b \geq +1 - \xi_i & \text{when } y_i = +1 \\ \omega^\top \mathbf{x}_i + b \leq -1 + \xi_i & \text{when } y_i = -1 \end{cases}$$

where $\xi_i \geq 0 \forall i$. There is a classification error when $\xi_i > 1$.

The idea is then to solve

$$\min \left\{ \frac{1}{2} \omega^\top \omega + C \mathbf{1}^\top \mathbf{1}_{\xi > 1} \right\}, \text{ instead of } \min \left\{ \frac{1}{2} \omega^\top \omega \right\}$$

Support Vector Machines, with a Linear Kernel

So far,

$$d(\mathbf{x}_0, H_{\omega, b}) = \min_{\mathbf{x} \in H_{\omega, b}} \{\|\mathbf{x}_0 - \mathbf{x}\|_{\ell_2}\}$$

where $\|\cdot\|_{\ell_2}$ is the Euclidean (ℓ_2) norm,

$$\|\mathbf{x}_0 - \mathbf{x}\|_{\ell_2} = \sqrt{(\mathbf{x}_0 - \mathbf{x}) \cdot (\mathbf{x}_0 - \mathbf{x})} = \sqrt{\mathbf{x}_0 \cdot \mathbf{x}_0 - 2\mathbf{x}_0 \cdot \mathbf{x} + \mathbf{x} \cdot \mathbf{x}}$$

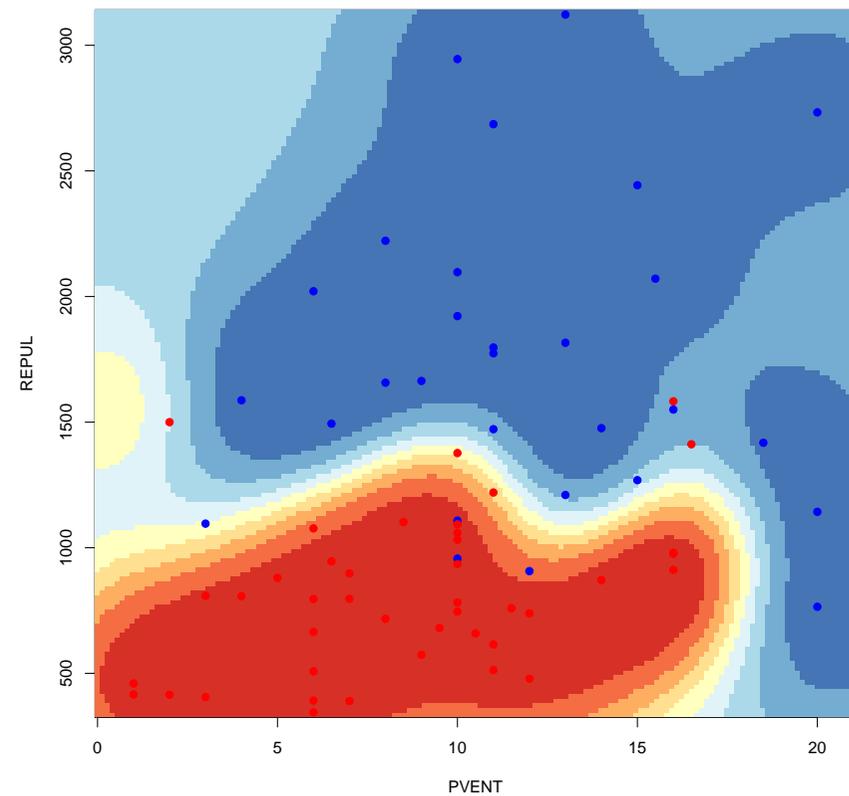
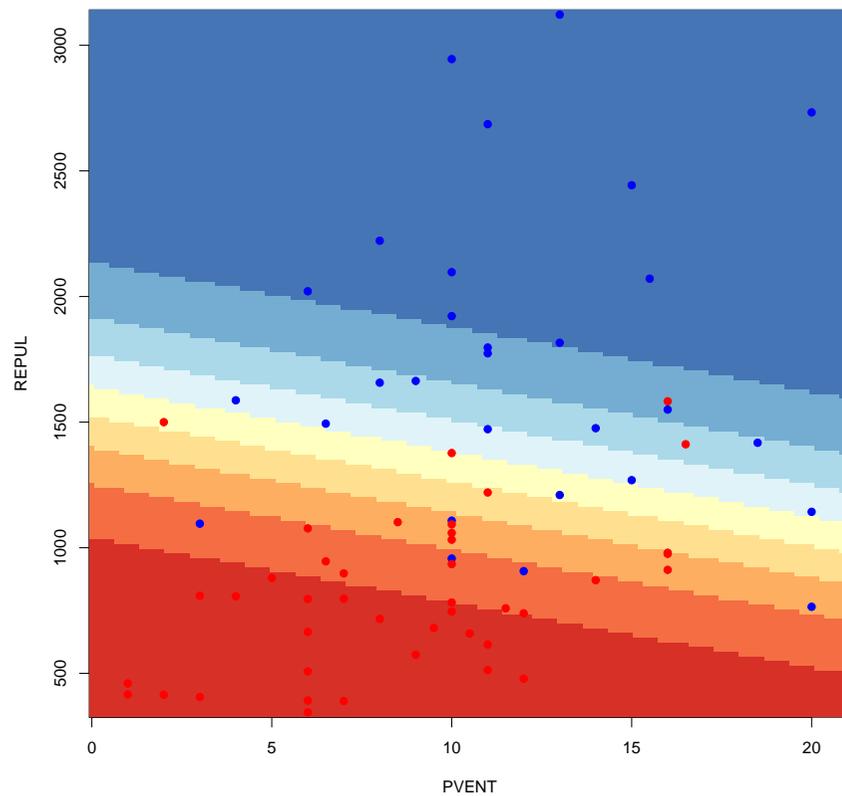
More generally,

$$d(\mathbf{x}_0, H_{\omega, b}) = \min_{\mathbf{x} \in H_{\omega, b}} \{\|\mathbf{x}_0 - \mathbf{x}\|_k\}$$

where $\|\cdot\|_k$ is some kernel-based norm,

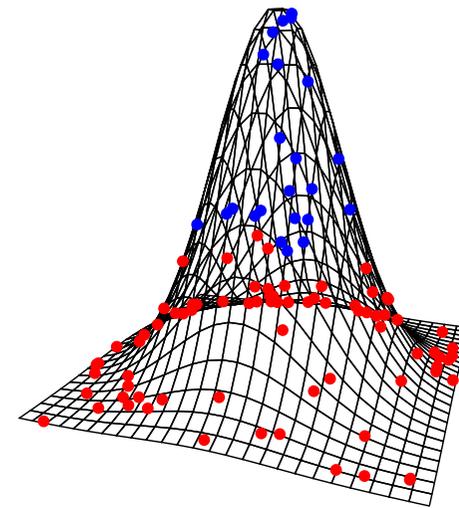
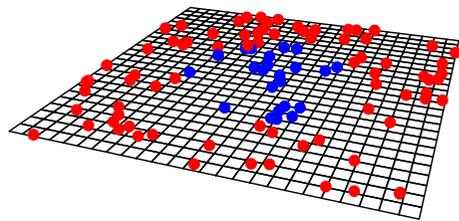
$$\|\mathbf{x}_0 - \mathbf{x}\|_k = \sqrt{k(\mathbf{x}_0, \mathbf{x}_0) - 2k(\mathbf{x}_0, \mathbf{x}) + k(\mathbf{x}, \mathbf{x})}$$

Support Vector Machines, with a Non Linear Kernel



Heuristics on SVMs

An interpretation is that data aren't linearly separable in the original space, but might be separable by some kernel transformation,



Penalization and Mean Square Error

Consider the quadratic loss function, $\ell(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2$, the risk function becomes the mean squared error of the estimate,

$$R(\theta, \hat{\theta}) = \mathbb{E}(\theta - \hat{\theta})^2 = \underbrace{[\theta - \mathbb{E}(\hat{\theta})]^2}_{\text{bias}^2} + \underbrace{\mathbb{E}(\mathbb{E}[\hat{\theta}] - \hat{\theta})^2}_{\text{variance}}$$

Get back to the initial example, $y_i \in \{0, 1\}$, with $p = \mathbb{P}(Y = 1)$.

Consider the estimate that minimizes the mse, that can be written $\hat{p} = (1 - \alpha)\bar{y}$, then

$$\text{mse}(\hat{p}) = \alpha^2 p^2 + (1 - \alpha)^2 \frac{p(1 - p)}{n}$$

$$\text{then } \alpha^* = \frac{1 - p}{1 + (n - 1)p}.$$

i.e. unbiased estimators have nice mathematical properties, but can be improved.

Linear Model

Consider some linear model $y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i$ for all $i = 1, \dots, n$.

Assume that ε_i are i.i.d. with $\mathbb{E}(\varepsilon) = 0$ (and finite variance). Write

$$\underbrace{\begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}}_{\mathbf{y}, n \times 1} = \underbrace{\begin{pmatrix} 1 & x_{1,1} & \cdots & x_{1,k} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n,1} & \cdots & x_{n,k} \end{pmatrix}}_{\mathbf{X}, n \times (k+1)} \underbrace{\begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}}_{\boldsymbol{\beta}, (k+1) \times 1} + \underbrace{\begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}}_{\boldsymbol{\varepsilon}, n \times 1}.$$

Assuming $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbb{I})$, the maximum likelihood estimator of $\boldsymbol{\beta}$ is

$$\hat{\boldsymbol{\beta}} = \operatorname{argmin}\{\|\mathbf{y} - \mathbf{X}^\top \boldsymbol{\beta}\|_{\ell_2}\} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

... under the assumption that $\mathbf{X}^\top \mathbf{X}$ is a full-rank matrix.

What if $\mathbf{X}_i^\top \mathbf{X}$ cannot be inverted? Then $\hat{\boldsymbol{\beta}} = [\mathbf{X}^\top \mathbf{X}]^{-1} \mathbf{X}^\top \mathbf{y}$ does not exist, but $\hat{\boldsymbol{\beta}}_\lambda = [\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I}]^{-1} \mathbf{X}^\top \mathbf{y}$ always exist if $\lambda > 0$.

Ridge Regression

The estimator $\hat{\beta} = [\mathbf{X}^T \mathbf{X} + \lambda \mathbb{I}]^{-1} \mathbf{X}^T \mathbf{y}$ is the **Ridge** estimate obtained as solution of

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^n [y_i - \beta_0 - \mathbf{x}_i^T \beta]^2 + \lambda \underbrace{\|\beta\|_{\ell_2}}_{\mathbf{1}^T \beta^2} \right\}$$

for some tuning parameter λ . One can also write

$$\hat{\beta} = \underset{\beta; \|\beta\|_{\ell_2} \leq s}{\operatorname{argmin}} \{ \|\mathbf{Y} - \mathbf{X}^T \beta\|_{\ell_2} \}$$

Remark Note that we solve $\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \{ \operatorname{objective}(\beta) \}$ where

$$\operatorname{objective}(\beta) = \underbrace{\mathcal{L}(\beta)}_{\text{training loss}} + \underbrace{\mathcal{R}(\beta)}_{\text{regularization}}$$

Going further on sparsity issues

In several applications, k can be (very) large, but a lot of features are just noise: $\beta_j = 0$ for many j 's. Let s denote the number of relevant features, with $s \ll k$, cf [Hastie, Tibshirani & Wainwright \(2015\)](#),

$$s = \text{card}\{\mathcal{S}\} \text{ where } \mathcal{S} = \{j; \beta_j \neq 0\}$$

The model is now $y = \mathbf{X}_{\mathcal{S}}^{\top} \boldsymbol{\beta}_{\mathcal{S}} + \varepsilon$, where $\mathbf{X}_{\mathcal{S}}^{\top} \mathbf{X}_{\mathcal{S}}$ is a full rank matrix.

Going further on sparsity issues

Define $\|\mathbf{a}\|_{\ell_0} = \sum \mathbf{1}(|a_i| > 0)$. Ici $\dim(\boldsymbol{\beta}) = s$.

We wish we could solve

$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta}; \|\boldsymbol{\beta}\|_{\ell_0} \leq s}{\operatorname{argmin}} \{ \|\mathbf{Y} - \mathbf{X}^\top \boldsymbol{\beta}\|_{\ell_2} \}$$

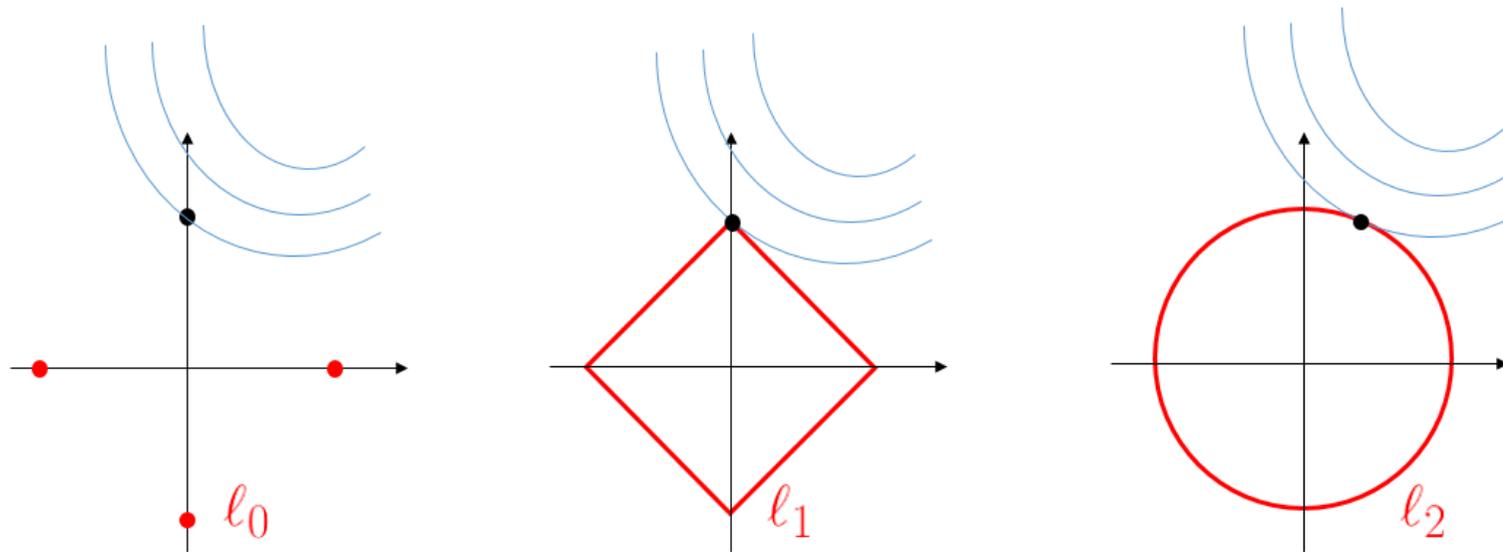
Problem: it is usually not possible to describe all possible constraints, since $\binom{s}{k}$ coefficients should be chosen here (with k (very) large).

Idea: solve the dual problem

$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta}; \|\mathbf{Y} - \mathbf{X}^\top \boldsymbol{\beta}\|_{\ell_2} \leq h}{\operatorname{argmin}} \{ \|\boldsymbol{\beta}\|_{\ell_0} \}$$

where we might convexify the ℓ_0 norm, $\|\cdot\|_{\ell_0}$.

Regularization l_0 , l_1 and l_2



$$\min\{\|\beta\|_{l_*}\} \text{ subject to } \|\mathbf{Y} - \mathbf{X}^T \beta\|_{l_2} \leq h$$

Going further on sparsity issues

On $[-1, +1]^k$, the convex hull of $\|\beta\|_{\ell_0}$ is $\|\beta\|_{\ell_1}$

On $[-a, +a]^k$, the convex hull of $\|\beta\|_{\ell_0}$ is $a^{-1}\|\beta\|_{\ell_1}$

Hence,

$$\hat{\beta} = \underset{\beta; \|\beta\|_{\ell_1} \leq \tilde{s}}{\operatorname{argmin}} \{ \|\mathbf{Y} - \mathbf{X}^T \beta\|_{\ell_2} \}$$

is equivalent (Kuhn-Tucker theorem) to the Lagrangian optimization problem

$$\hat{\beta} = \operatorname{argmin} \{ \|\mathbf{Y} - \mathbf{X}^T \beta\|_{\ell_2} + \lambda \|\beta\|_{\ell_1} \}$$

LASSO *Least Absolute Shrinkage and Selection Operator*

$$\hat{\beta} \in \operatorname{argmin}\{\|\mathbf{Y} - \mathbf{X}^T \beta\|_{\ell_2} + \lambda \|\beta\|_{\ell_1}\}$$

is a convex problem (several algorithms^{*}), but not strictly convex (no unicity of the minimum). Nevertheless, predictions $\hat{\mathbf{y}} = \mathbf{x}^T \hat{\beta}$ are unique

^{*} MM, minimize majorization, coordinate descent [Hunter \(2003\)](#).

Optimal LASSO Penalty

Use cross validation, e.g. K -fold,

$$\hat{\beta}_{(-k)}(\lambda) = \operatorname{argmin} \left\{ \sum_{i \notin \mathcal{I}_k} [y_i - \mathbf{x}_i^\top \beta]^2 + \lambda \|\beta\| \right\}$$

then compute the sum of the squared errors,

$$Q_k(\lambda) = \sum_{i \in \mathcal{I}_k} [y_i - \mathbf{x}_i^\top \hat{\beta}_{(-k)}(\lambda)]^2$$

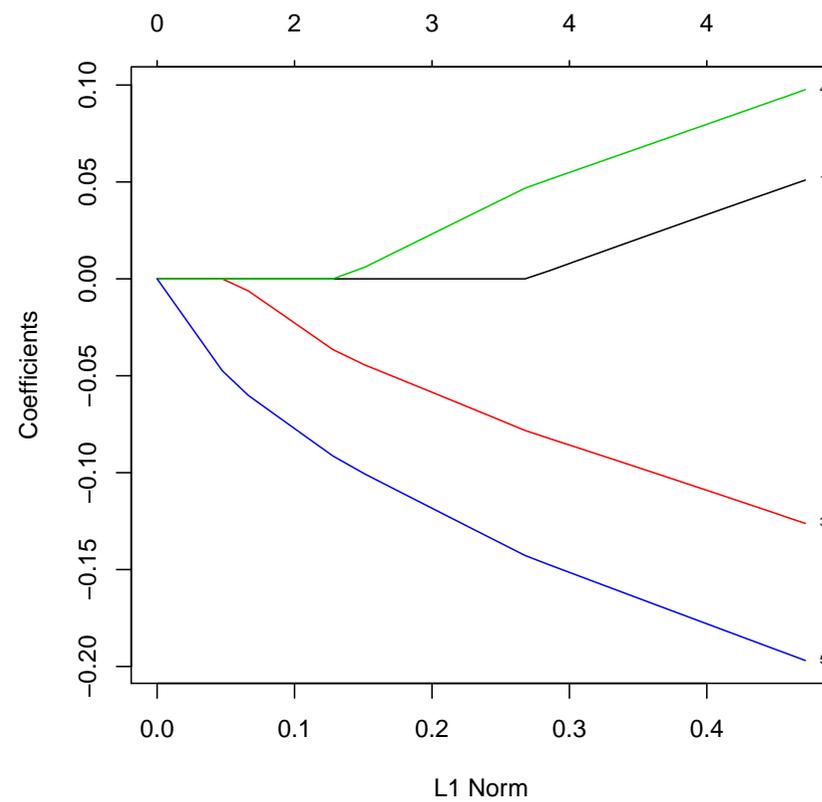
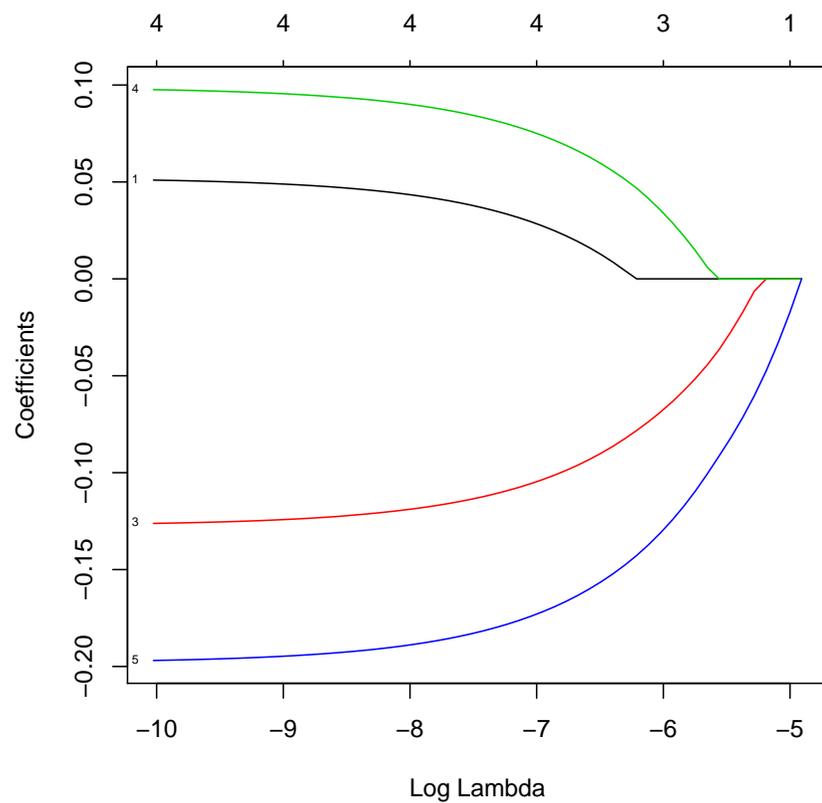
and finally solve

$$\lambda^* = \operatorname{argmin} \left\{ \bar{Q}(\lambda) = \frac{1}{K} \sum_k Q_k(\lambda) \right\}$$

Note that this might overfit, so [Hastie, Tibshiriani & Friedman \(2009\)](#) suggest the largest λ such that

$$\bar{Q}(\lambda) \leq \bar{Q}(\lambda^*) + \operatorname{se}[\lambda^*] \quad \text{with} \quad \operatorname{se}[\lambda]^2 = \frac{1}{K^2} \sum_{k=1}^K [Q_k(\lambda) - \bar{Q}(\lambda)]^2$$

LASSO



Penalization and GLM's

The logistic regression is based on empirical risk, when $y \in \{0, 1\}$

$$-\frac{1}{n} \sum_{i=1}^n (y_i \mathbf{x}_i^T \boldsymbol{\beta} - \log[1 + \exp(\mathbf{x}_i^T \boldsymbol{\beta})])$$

or, if $y \in \{-1, +1\}$,

$$\frac{1}{n} \sum_{i=1}^n \log [1 + \exp(y_i \mathbf{x}_i^T \boldsymbol{\beta})] .$$

A regularized version with the ℓ_1 norm is the **LASSO** logistic regression

$$\frac{1}{n} \sum_{i=1}^n \log [1 + \exp(y_i \mathbf{x}_i^T \boldsymbol{\beta})] + \lambda \|\boldsymbol{\beta}\|_1$$

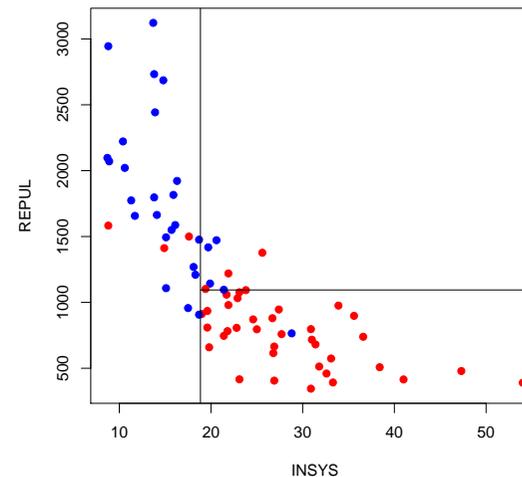
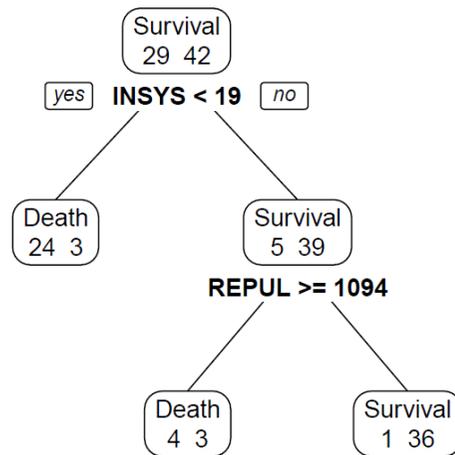
or more generally, with smoothing functions

$$\frac{1}{n} \sum_{i=1}^n \log [1 + \exp(y_i g(\mathbf{x}_i))] + \lambda \|g\|$$

Classification (and Regression) Trees, CART

one of the predictive modelling approaches used in statistics, data mining and machine learning [...] In tree structures, leaves represent class labels and branches represent conjunctions of features that lead to those class labels.

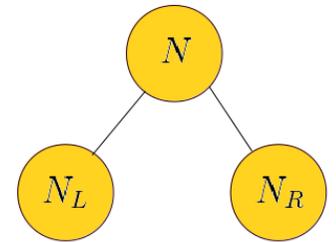
(Source: [wikipedia](https://en.wikipedia.org/wiki/Decision_tree_learning)).



Classification (and Regression) Trees, CART

To split N into two $\{N_L, N_R\}$, consider

$$\mathcal{I}(N_L, N_R) = \sum_{x \in \{L, R\}} \frac{n_x}{n} \mathcal{I}(N_x)$$



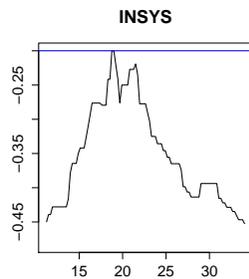
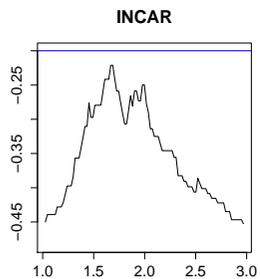
e.g. Gini index (used originally in CART, see [Breiman et al. \(1984\)](#))

$$\text{gini}(N_L, N_R) = - \sum_{x \in \{L, R\}} \frac{n_x}{n} \sum_{y \in \{0, 1\}} \frac{n_{x,y}}{n_x} \left(1 - \frac{n_{x,y}}{n_x} \right)$$

and the cross-entropy (used in C4.5 and C5.0)

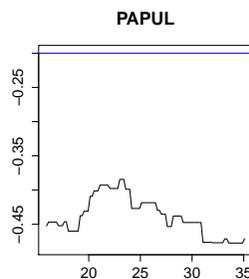
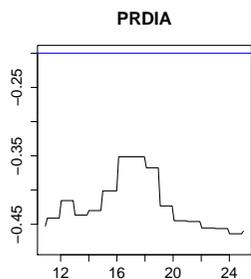
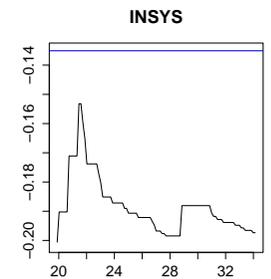
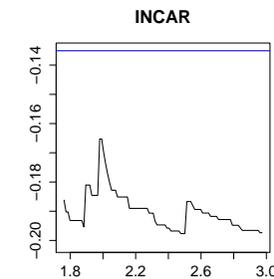
$$\text{entropy}(N_L, N_R) = - \sum_{x \in \{L, R\}} \frac{n_x}{n} \sum_{y \in \{0, 1\}} \frac{n_{x,y}}{n_x} \log \left(\frac{n_{x,y}}{n_x} \right)$$

Classification (and Regression) Trees, CART

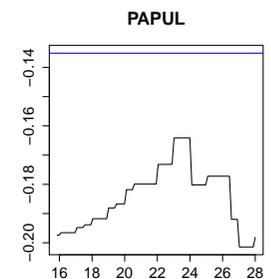
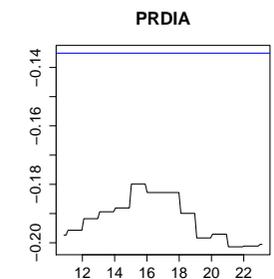


$$N_L: \{x_{i,j} \leq s\} \quad N_R: \{x_{i,j} > s\}$$

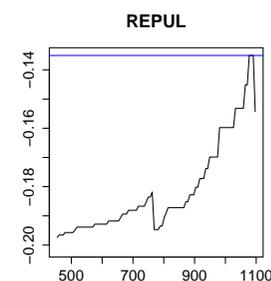
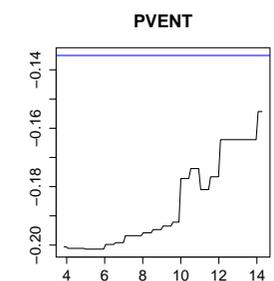
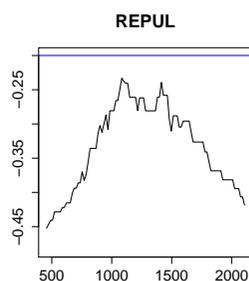
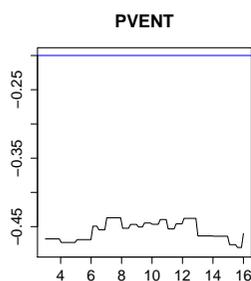
$$\text{solve } \max_{j \in \{1, \dots, k\}, s} \{\mathcal{I}(N_L, N_R)\}$$



← first split



second split →



Pruning Trees

One can grow a big tree, until leaves have a (preset) small number of observations, and then possibly go back and prune branches (or leaves) that do not improve gains on good classification sufficiently.

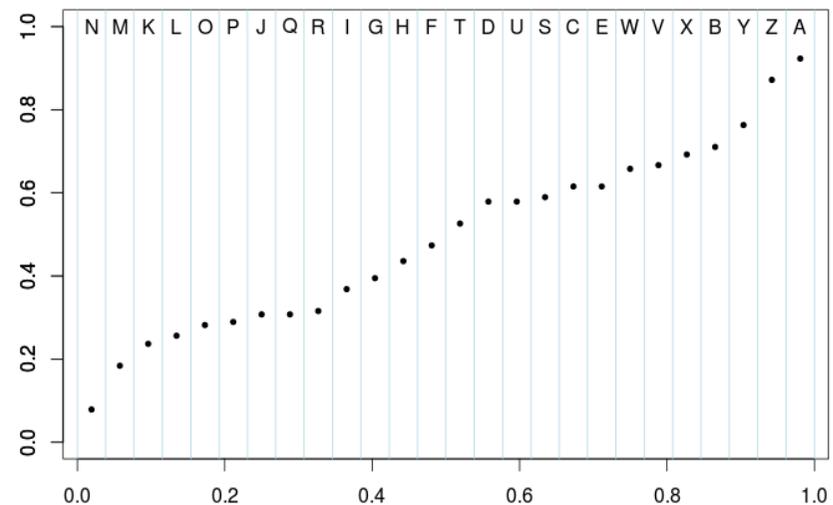
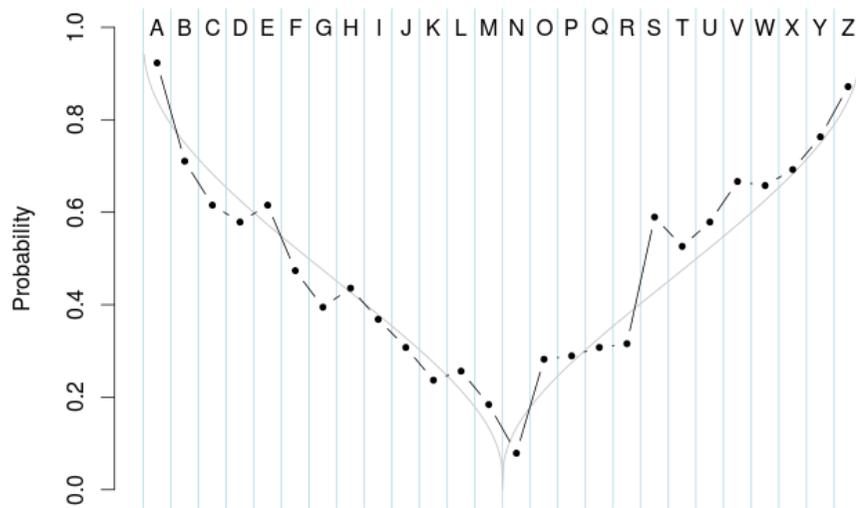
Or we can decide, at each node, whether we split, or not.

In trees, overfitting increases with the number of steps, and leaves. Drop in impurity at node N is defined as

$$\Delta\mathcal{I}(N_L, N_R) = \mathcal{I}(N) - \mathcal{I}(N_L, N_R) = \mathcal{I}(N) - \left(\frac{n_L}{n} \mathcal{I}(N_L) - \frac{n_R}{n} \mathcal{I}(N_R) \right)$$

(Fast) Trees with Categorical Features

Consider some simple categorical covariate, $x \in \{A, B, C, \dots, Y, Z\}$, defined from a continuous latent variable $\tilde{x} \sim \mathcal{U}([0, 1])$.



Compute $\bar{y}(x) = \frac{1}{n_x} \sum_{i:x_i=x} y_i \approx \mathbb{E}[Y|X = x]$ and sort them

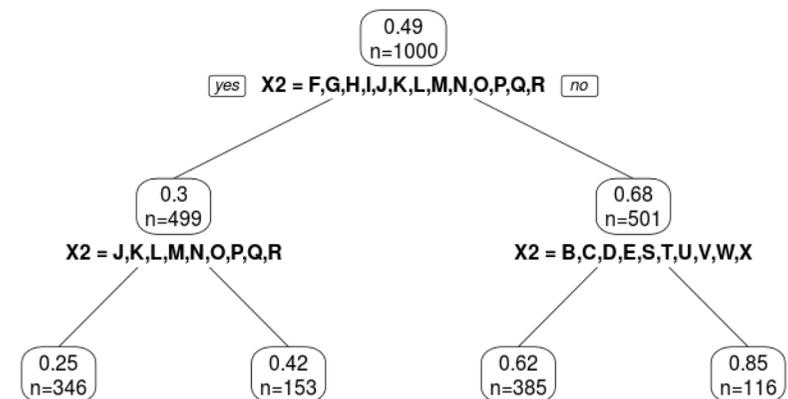
$$\bar{y}(x_{1:26}) \leq \bar{y}(x_{2:26}) \leq \dots \leq \bar{y}(x_{25:26}) \leq \bar{y}(x_{26:26}).$$

(Fast) Trees with Categorical Features

Then the split is done base on sample

$$x \in \{x_{1:26}, \dots, x_{j:26}\}$$

$$\text{vs. } x \in \{x_{j+1:26}, \dots, x_{26:26}\}$$



Bagging

Bootstrapped Aggregation (Bagging) , is a machine learning ensemble meta-algorithm designed to improve the stability and accuracy of machine learning algorithms used in statistical classification (Source: [wikipedia](#)).

It is an ensemble method that creates multiple models of the same type from different sub-samples of the same dataset [**bootstrap**]. The predictions from each separate model are combined together to provide a superior result [**aggregation**].

→ can be used on any kind of model, but interesting for trees, see [Breiman \(1996\)](#)

Bootstrap can be used to define the concept of **margin**,

$$\text{margin}_i = \frac{1}{B} \sum_{b=1}^B \mathbf{1}(\hat{y}_i = y_i) - \frac{1}{B} \sum_{b=1}^B \mathbf{1}(\hat{y}_i \neq y_i)$$

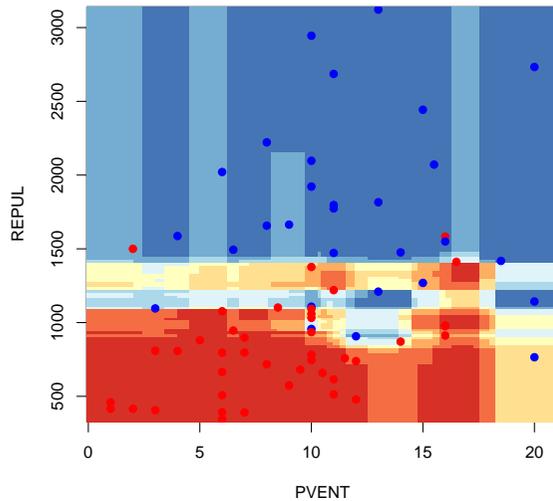
Remark Probability that i th row is not selection $(1 - n^{-1})^n \rightarrow e^{-1} \sim 36.8\%$, cf training / validation samples (2/3-1/3)

Bagging : Bootstrap Aggregation

For classes, $\tilde{m}(\mathbf{x}) = \operatorname{argmax}_y \sum_{b=1}^B \mathbf{1}(y = \hat{m}^{(b)})$.

For probabilities,

$$\tilde{m}(\mathbf{x}) = \frac{1}{n} \sum_{b=1}^B \hat{m}^{(b)}(\mathbf{x}) = \frac{1}{n} \sum_{b=1}^B \sum_{j=1}^{k_b} y_j \mathbf{1}(\mathbf{x}_i \in C_j).$$



Model Selection and Gini/Lorentz (on incomes)

Consider an ordered sample $\{y_1, \dots, y_n\}$, then Lorenz curve is

$$\{F_i, L_i\} \text{ with } F_i = \frac{i}{n} \text{ and } L_i = \frac{\sum_{j=1}^i y_j}{\sum_{j=1}^n y_j}$$

The theoretical curve, given a distribution F , is

$$u \mapsto L(u) = \frac{\int_{-\infty}^{F^{-1}(u)} t dF(t)}{\int_{-\infty}^{+\infty} t dF(t)}$$

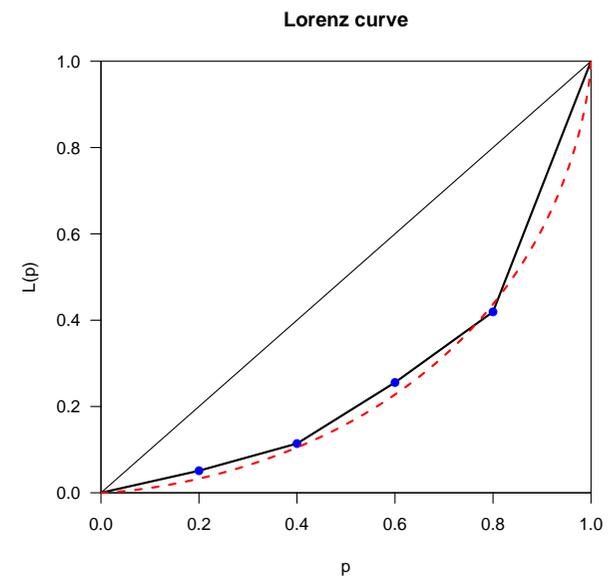
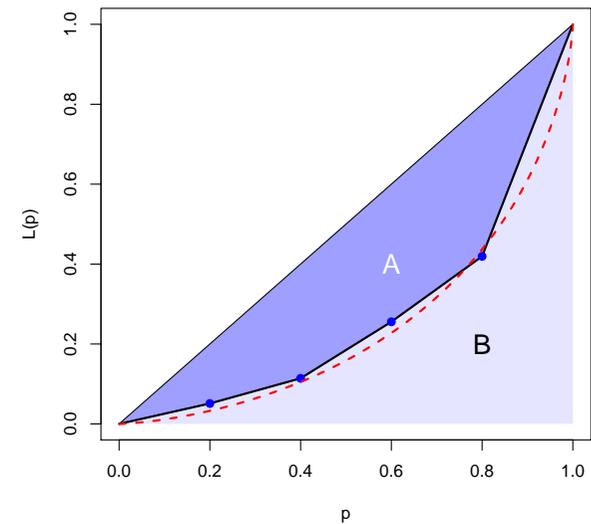
see [Gastwirth \(1972\)](#)

Model Selection and Gini/Lorentz

Gini index is the ratio of the areas $\frac{A}{A+B}$. Thus,

$$G = \frac{2}{n(n-1)\bar{x}} \sum_{i=1}^n i \cdot x_{i:n} - \frac{n+1}{n-1}$$

$$= \frac{1}{\mathbb{E}(Y)} \int_0^{\infty} F(y)(1-F(y))dy$$



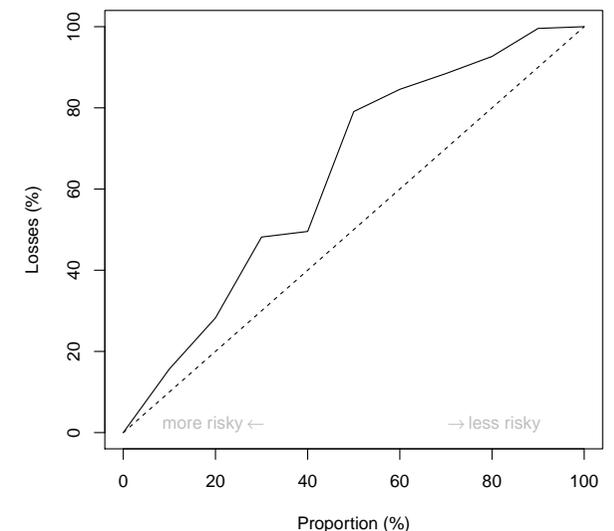
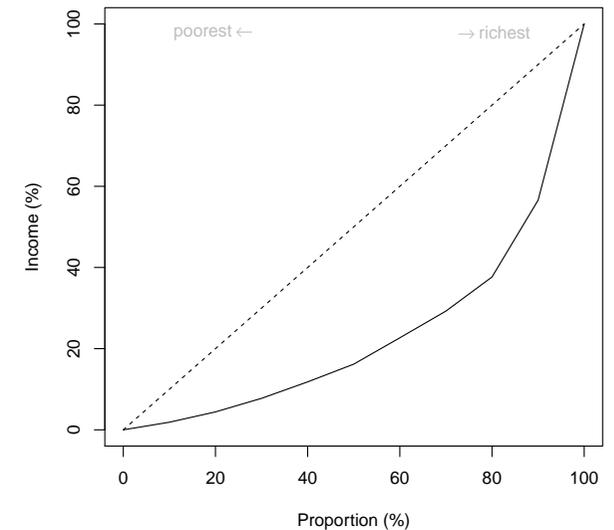
Model Selection

Consider an ordered sample $\{y_1, \dots, y_n\}$ of incomes, with $y_1 \leq y_2 \leq \dots \leq y_n$, then Lorenz curve is

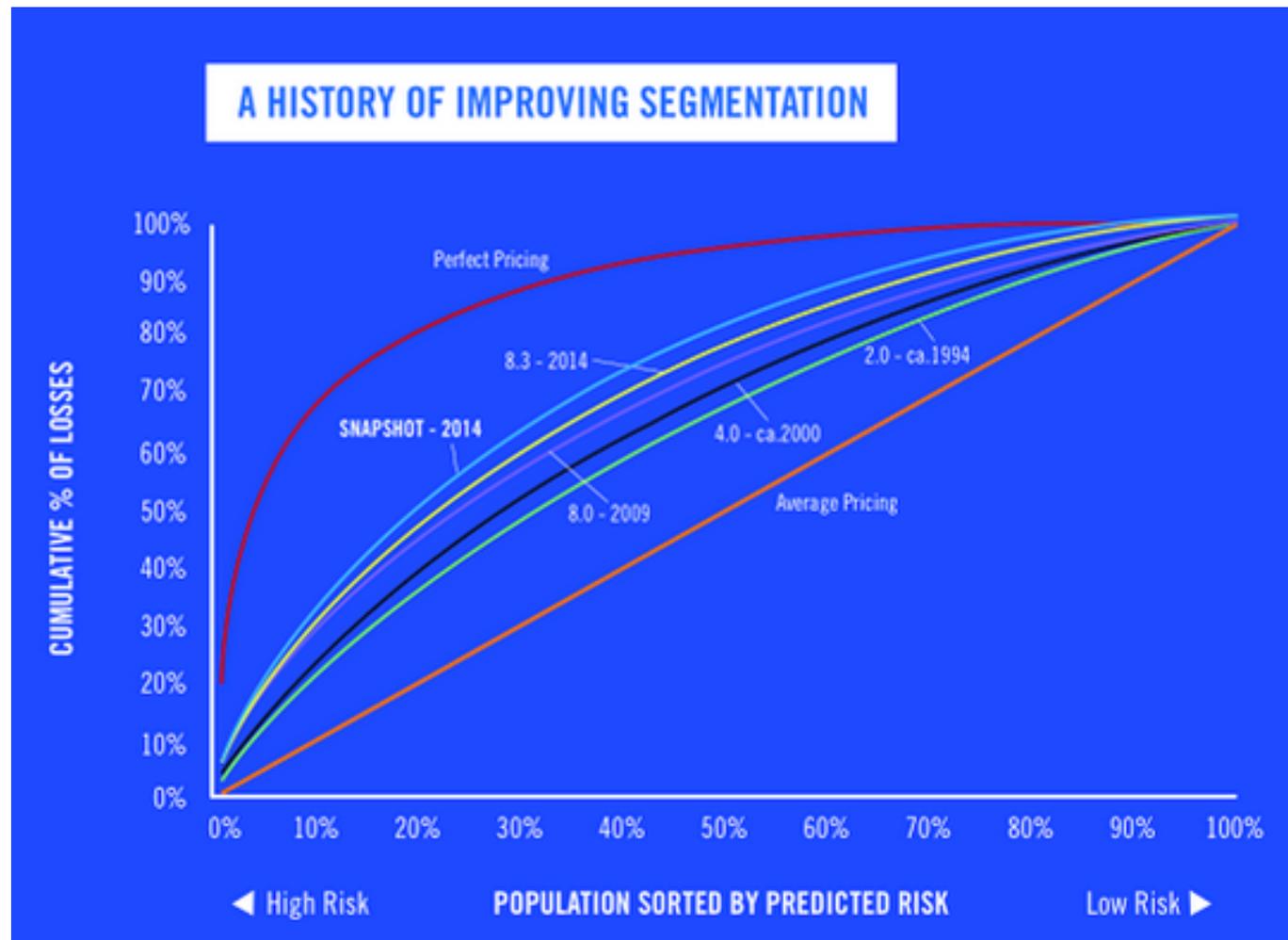
$$\{F_i, L_i\} \text{ with } F_i = \frac{i}{n} \text{ and } L_i = \frac{\sum_{j=1}^i y_j}{\sum_{j=1}^n y_j}$$

We have observed losses y_i and premiums $\hat{\pi}(x_i)$. Consider an **ordered sample by the model**, see [Frees, Meyers & Cummins \(2014\)](#), $\hat{\pi}(x_1) \geq \hat{\pi}(x_2) \geq \dots \geq \hat{\pi}(x_n)$, then plot

$$\{F_i, L_i\} \text{ with } F_i = \frac{i}{n} \text{ and } L_i = \frac{\sum_{j=1}^i y_j}{\sum_{j=1}^n y_j}$$



Model Selection



See Frees *et al.* (2010) or Tevet (2013).

Part 4. Small Data and Bayesian Philosophy

“it’s time to adopt modern Bayesian data analysis as standard procedure in our scientific practice and in our educational curriculum. Three reasons:

1. Scientific disciplines from astronomy to zoology are moving to Bayesian analysis.
We should be leaders of the move, not followers.
2. Modern Bayesian methods provide richer information, with greater flexibility and broader applicability than 20th century methods. Bayesian methods are intellectually coherent and intuitive.

Bayesian analyses are readily computed with modern software and hardware.

3. Null-hypothesis significance testing (NHST), with its reliance on p values, has many problems.

There is little reason to persist with NHST now that Bayesian methods are accessible to everyone.

My conclusion from those points is that we should do whatever we can to encourage the move to Bayesian data analysis.” John Kruschke,

(quoted in Meyers & Guszczka (2013))

Bayes vs. Frequentist, inference on heads/tails

Consider some Bernoulli sample $\mathbf{x} = \{x_1, x_2, \dots, x_n\}$, where $x_i \in \{0, 1\}$.

X_i 's are i.i.d. $\mathcal{B}(p)$ variables, $f_X(x) = p^x[1-p]^{1-x}$, $x \in \{0, 1\}$.

Standard frequentist approach

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n x_i = \operatorname{argmax}_{p \in (0,1)} \left\{ \underbrace{\prod_{i=1}^n f_X(x_i)}_{\mathcal{L}(p; \mathbf{x})} \right\}$$

From the central limit theorem

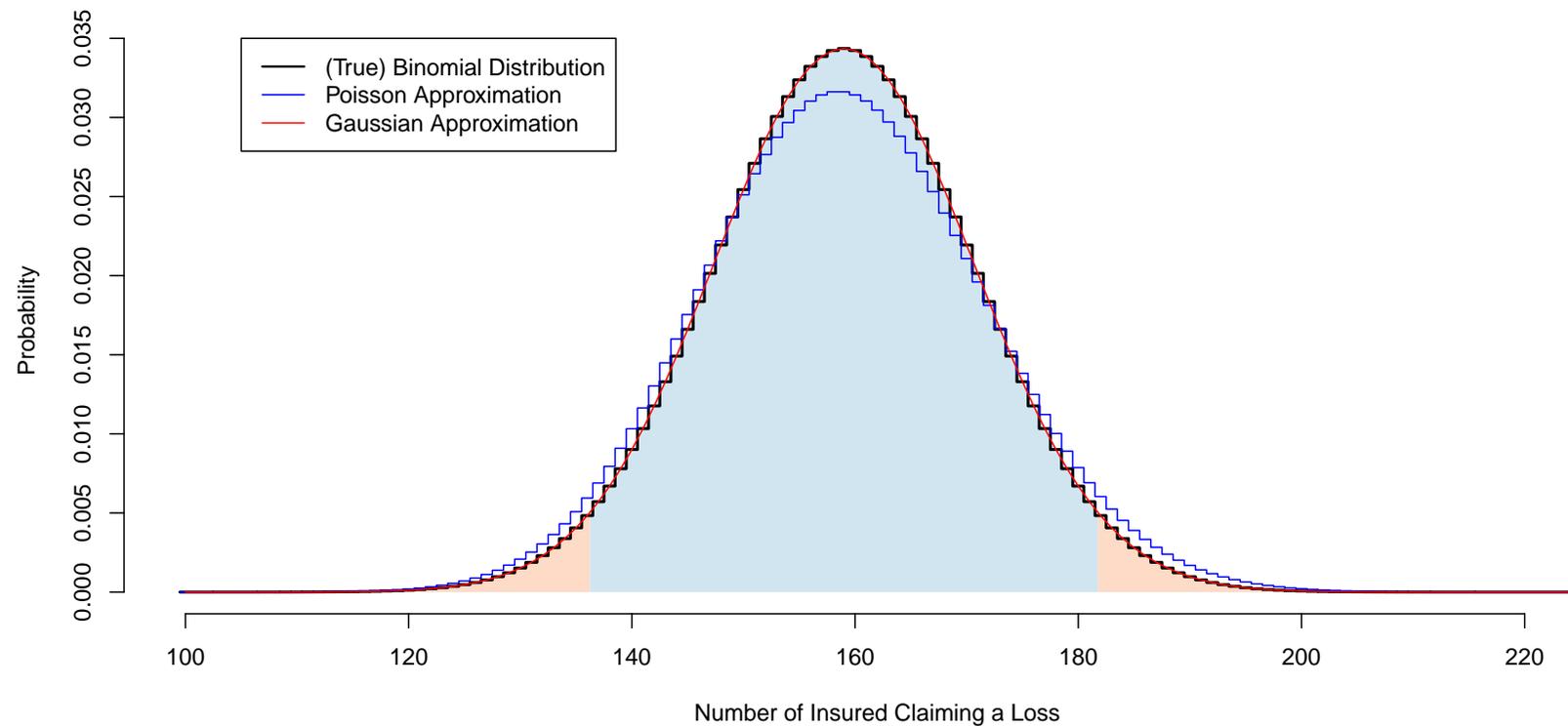
$$\sqrt{n} \frac{\hat{p} - p}{\sqrt{p(1-p)}} \xrightarrow{\mathcal{L}} \mathcal{N}(0, 1) \text{ as } n \rightarrow \infty$$

we can derive an approximated 95% confidence interval

$$\left[\hat{p} \pm \frac{1.96}{\sqrt{n}} \sqrt{\hat{p}(1-\hat{p})} \right]$$

Bayes vs. Frequentist, inference on heads/tails

Example out of 1,047 contracts, 159 claimed a loss



Small Data and Black Swans

Example [Operational risk] What if our sample is $\mathbf{x} = \{0, 0, 0, 0, 0\}$?

How would we derive a confidence interval for p ?

“INA’s chief executive officer, dressed as Santa Claus, asked an unthinkable question: Could anyone predict the probability of two planes colliding in midair? Santa was asking his chief actuary, L. H. Longley-Cook, to make a prediction based on no experience at all. There had never been a serious midair collision of commercial planes. Without any past experience or repetitive experimentation, any orthodox statistician had to answer Santa’s question with a resounding no.”

the theory 
 that would
 not die 
 how bayes' rule cracked
 the enigma code,
 hunted down russian
 submarines & emerged
 triumphant from two 
 centuries of controversy
 sharon bertsch mcgrayne

Bayes, the theory that would not die

Liu et al. (1996) claim that “ Statistical methods with a Bayesian flavor [...] have long been used in the insurance industry”.

History of Bayesian statistics, *the theory that would not die* by Sharon Bertsch McGrayne

“[Arthur] Bailey spent his first year in New York [in 1918] trying to prove to himself that ‘all of the fancy actuarial [Bayesian] procedures of the casualty business were mathematically unsound.’ After a year of intense mental struggle, however, realized to his consternation that actuarial sledgehammering worked” [...]

Bayes, the theory that would not die

[...] “ He even preferred it to the elegance of frequentism. He positively liked formulae that described ‘actual data . . . I realized that the hard-shelled underwriters were recognizing certain facts of life neglected by the statistical theorists.’ He wanted to give more weight to a large volume of data than to the frequentists small sample; doing so felt surprisingly ‘logical and reasonable’. He concluded that only a ‘suicidal’ actuary would use Fishers method of maximum likelihood, which assigned a zero probability to nonevents. Since many businesses file no insurance claims at all, Fishers method would produce premiums too low to cover future losses.”

Bayes's theorem

Consider some hypothesis H and some evidence E , then

$$\mathbb{P}_E(H) = \mathbb{P}(H|E) = \frac{\mathbb{P}(H \cap E)}{\mathbb{P}(E)} = \frac{\mathbb{P}(H) \cdot \mathbb{P}(E|H)}{\mathbb{P}(E)}$$

Bayes rule,

$$\left\{ \begin{array}{l} \text{prior probability } \mathbb{P}(H) \\ \text{versus posterior probability after receiving evidence } E, \mathbb{P}_E(H) = \mathbb{P}(H|E). \end{array} \right.$$

In Bayesian (parametric) statistics, $H = \{\theta \in \Theta\}$ and $E = \{\mathbf{X} = \mathbf{x}\}$.

Bayes' Theorem,

$$\pi(\theta|\mathbf{x}) = \frac{\pi(\theta) \cdot f(\mathbf{x}|\theta)}{f(\mathbf{x})} = \frac{\pi(\theta) \cdot f(\mathbf{x}|\theta)}{\int f(\mathbf{x}|\theta)\pi(\theta)d\theta} \propto \pi(\theta) \cdot f(\mathbf{x}|\theta)$$

Small Data and Black Swans

Consider sample $\mathbf{x} = \{0, 0, 0, 0, 0\}$.

Here the likelihood is

$$\begin{cases} f(x_i|\theta) = \theta^{x_i} [1 - \theta]^{1-x_i} \\ f(\mathbf{x}|\theta) = \theta^{\mathbf{x}^\top \mathbf{1}} [1 - \theta]^{n - \mathbf{x}^\top \mathbf{1}} \end{cases}$$

and we need a priori distribution $\pi(\cdot)$ e.g.
a beta distribution

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

$$\pi(\theta|\mathbf{x}) = \frac{\theta^{\alpha + \mathbf{x}^\top \mathbf{1}} [1 - \theta]^{\beta + n - \mathbf{x}^\top \mathbf{1}}}{B(\alpha + \mathbf{x}^\top \mathbf{1}, \beta + n - \mathbf{x}^\top \mathbf{1})}$$

On Bayesian Philosophy, Confidence vs. Credibility

for frequentists, a probability is a measure of the the frequency of repeated events

→ parameters are fixed (but unknown), and data are random

for Bayesians, a probability is a measure of the degree of certainty about values

→ parameters are random and data are fixed

“Bayesians : Given our observed data, there is a 95% probability that the true value of θ falls within the credible region

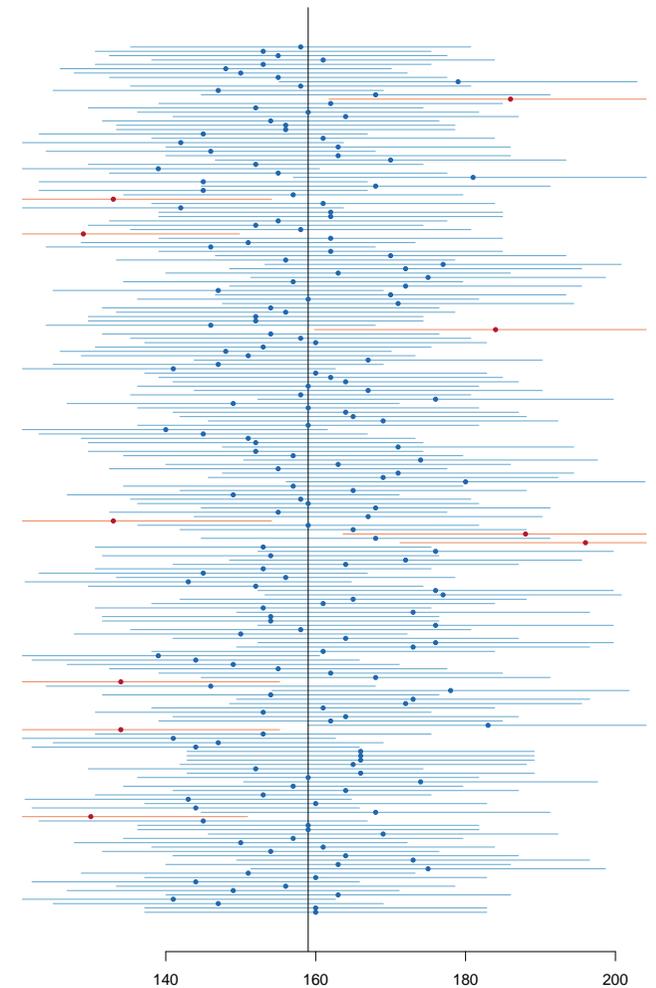
vs. Frequentists : There is a 95% probability that when I compute a confidence interval from data of this sort, the true value of θ will fall within it.” in [Vanderplas \(2014\)](#)

Example see [Jaynes \(1976\)](#), e.g. the truncated exponential

On Bayesian Philosophy, Confidence vs. Credibility

Example What is a 95% confidence interval of a proportion ? Here $\bar{x} = 159$ and $n = 1047$.

1. draw sets $(\tilde{x}_1, \dots, \tilde{x}_n)_k$ with $X_i \sim \mathcal{B}(\bar{x}/n)$
 2. compute for each set of values confidence intervals
 3. determine the fraction of these confidence interval that contain \bar{x}
- the parameter is fixed, and we guarantee that 95% of the confidence intervals will contain it.



On Bayesian Philosophy, Confidence vs. Credibility

Example What is 95% credible region of a proportion ? Here $\bar{x} = 159$ and $n = 1047$.

1. draw random parameters p_k with from the posterior distribution, $\pi(\cdot|\mathbf{x})$
 2. sample sets $(\tilde{x}_1, \dots, \tilde{x}_n)_k$ with $X_{i,k} \sim \mathcal{B}(p_k)$
 3. compute for each set of values means \bar{x}_k
 4. look at the proportion of those \bar{x}_k that are within this credible region $[\Pi^{-1}(.025|\mathbf{x}); \Pi^{-1}(.975|\mathbf{x})]$
- the credible region is fixed, and we guarantee that 95% of possible values of \bar{x} will fall within it.

Difficult concepts ? Difficult computations ?

We have a sample $\mathbf{x} = \{x_1, \dots, x_n\}$ i.i.d. from distribution $f_\theta(\cdot)$.

In predictive modeling, we need $\mathbb{E}(g(X)|\mathbf{x}) = \int g(x) f_{\theta|\mathbf{x}}(x) dx$ where

$$f_{\theta|\mathbf{x}}(x) = \int f_\theta(x) \cdot \pi(\theta|\mathbf{x}) d\theta$$

while prior density (without information \mathbf{x}) was

$$f_\theta(x) = \int f_\theta(x) \cdot \pi(\theta) d\theta$$

How can we derive $\pi(\theta|\mathbf{x})$?

Can we sample from $\pi(\theta|\mathbf{x})$ (use monte carlo technique to approximate the integral) ?

Computations not that simple... until the 90's : **MCMC**

Markov Chain

Stochastic process, $(X_t)_{t \in \mathbb{N}_*}$, on some discrete space Ω

$$\mathbb{P}(X_{t+1} = y | X_t = x, \underline{\mathbf{X}}_{t-1} = \underline{\mathbf{x}}_{t-1}) = \mathbb{P}(X_{t+1} = y | X_t = x) = P(x, y)$$

where P is a transition probability, that can be stored in a transition matrix, $\mathbf{P} = [P_{x,y}] = [P(x, y)]$.

Observe that $\mathbb{P}(X_{t+k} = y | X_t = x) = P_k(x, y)$ where $\mathbf{P}^k = [P_k(x, y)]$.

Under some condition, $\lim_{n \rightarrow \infty} \mathbf{P}^n = \mathbf{\Lambda} = [\boldsymbol{\lambda}^\top]$,

Problem given a distribution $\boldsymbol{\lambda}$, is it possible to generate a Markov Chain that converges to this distribution ?

Bonus Malus and Markov Chains

Ex no-claim bonus, see [Lemaire \(1995\)](#).

HONG KONG

Table B-9. Hong Kong System

Class	Premium	Class After		
		0	1 Claims	≥ 2
6	100	5	6	6
5	80	4	6	6
4	70	3	6	6
3	60	2	6	6
2	50	1	4	6
1	40	1	3	6

Starting class: 6.

Assume that the number of claims is $N \sim \mathcal{P}(21.7\%)$, so that $\mathbb{P}(N = 0) = 80\%$.

Hastings-Metropolis

Back to our problem, we want to sample from $\pi(\theta|\mathbf{x})$

i.e. generate $\theta_1, \dots, \theta_n, \dots$ from $\pi(\theta|\mathbf{x})$.

Hastings-Metropolis sampler will generate a Markov Chain (θ_t) as follows,

- generate θ_1
- generate θ^* and $U \sim \mathcal{U}([0, 1])$,

$$\text{compute } R = \frac{\pi(\theta^*|\mathbf{x})}{\pi(\theta_t|\mathbf{x})} \frac{P(\theta_t|\theta^*)}{P(\theta^*|\theta_{t-1})}$$

if $U < R$ set $\theta_{t+1} = \theta^*$

if $U \geq R$ set $\theta_{t+1} = \theta_t$

R is the acceptance ratio, we accept the new state θ^* with probability $\min\{1, R\}$.

Hastings-Metropolis

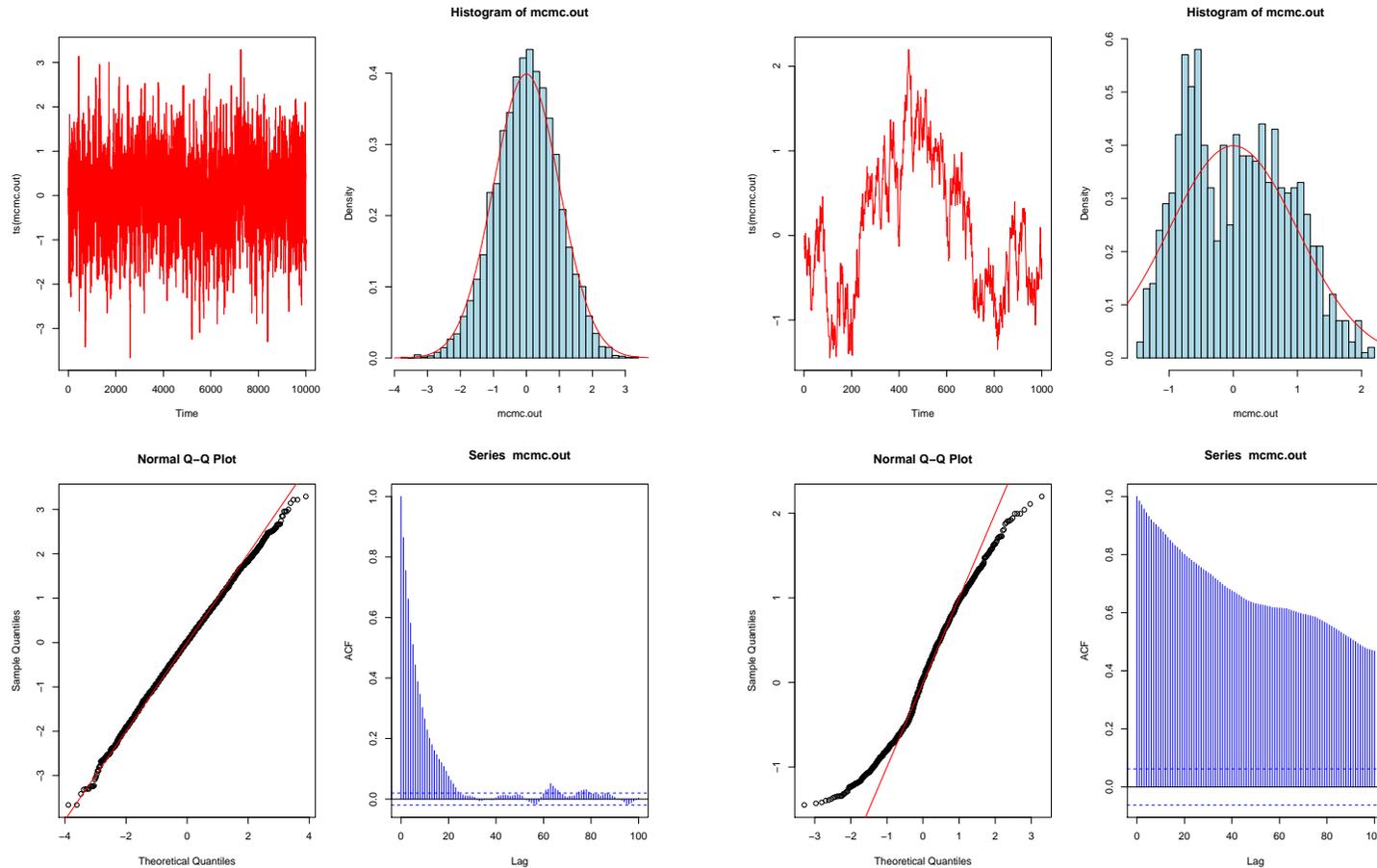
Observe that

$$R = \frac{\pi(\theta^*) \cdot f(\mathbf{x}|\theta^*)}{\pi(\theta_t) \cdot f(\mathbf{x}|\theta_t)} \frac{P(\theta_t|\theta^*)}{P(\theta^*|\theta_{t-1})}$$

In a more general case, we can have a Markov process, not a Markov chain.

E.g. $P(\theta^*|\theta_t) \sim \mathcal{N}(\theta_t, 1)$

Using MCMC to generate Gaussian values



Heuristics on Hastings-Metropolis

In standard Monte Carlo, generate θ_i 's i.i.d., then

$$\frac{1}{n} \sum_{i=1}^n g(\theta_i) \rightarrow \mathbb{E}[g(\theta)] = \int g(\theta) \pi(\theta) d\theta$$

(strong law of large numbers).

Well-behaved Markov Chains (\mathbf{P} aperiodic, irreducible, positive recurrent) can satisfy some ergodic property, similar to that LLN. More precisely,

- \mathbf{P} has a unique stationary distribution λ , i.e. $\lambda = \lambda \times \mathbf{P}$
- ergodic theorem

$$\frac{1}{n} \sum_{i=1}^n g(\theta_i) \rightarrow \int g(\theta) \lambda(\theta) d\theta$$

even if θ_i 's are not independent.

Heuristics on Hastings-Metropolis

Remark The conditions mentioned above are

- aperiodic, the chain does not regularly return to any state in multiples of some k .
- irreducible, the state can go from any state to any other state in some finite number of steps
- positively recurrent, the chain will return to any particular state with probability 1, and finite expected return time

Gibbs Sampler

For a multivariate problem, it is possible to use Gibbs sampler.

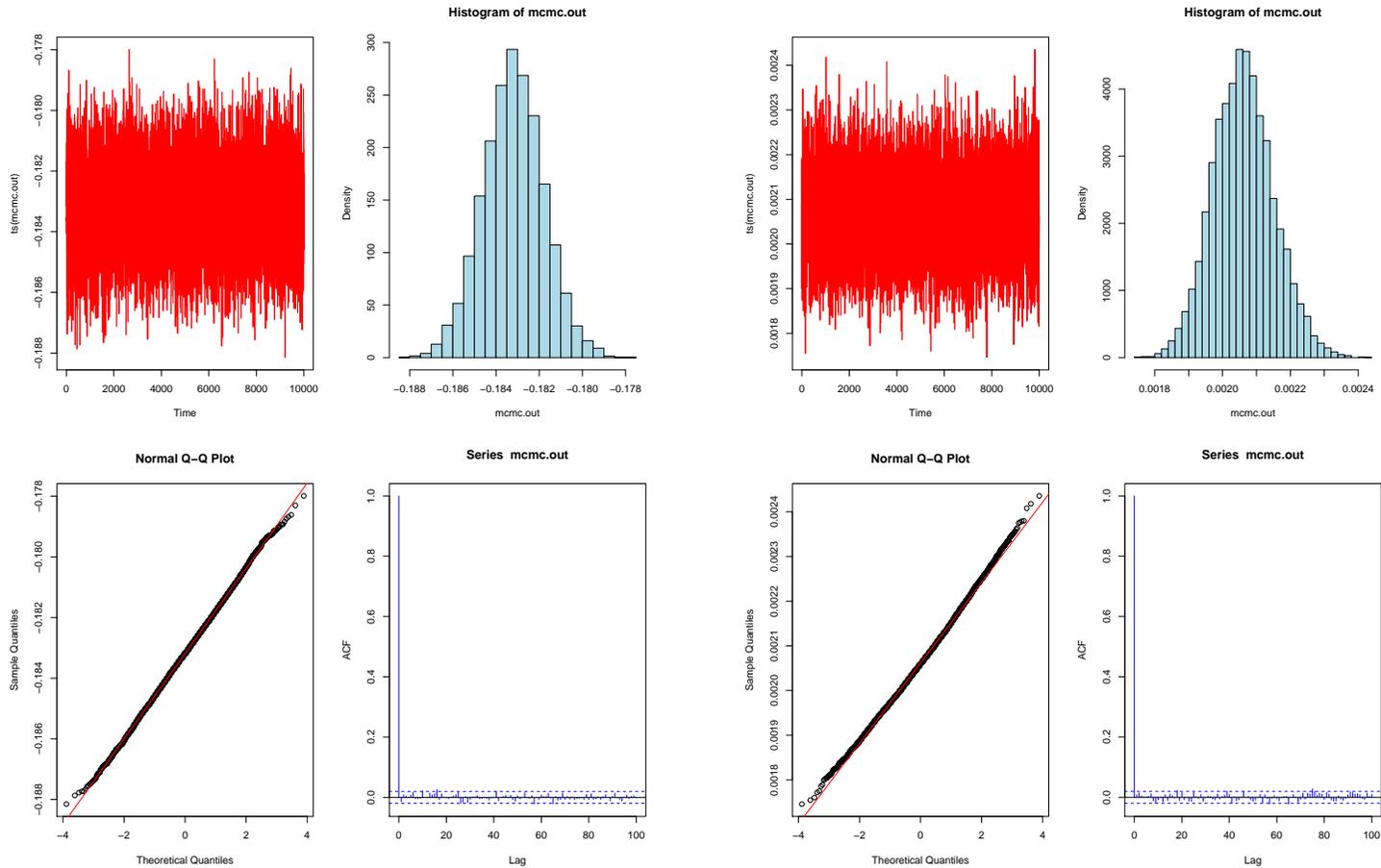
Example Assume that the loss ratio of a company has a lognormal distribution, $LN(\mu, \sigma^2)$, .e.g

Example Assume that we have a sample \mathbf{x} from a $\mathcal{N}(\mu, \sigma^2)$. We want the posterior distribution of $\boldsymbol{\theta} = (\mu, \sigma^2)$ given \mathbf{x} . Observe here that if priors are Gaussian $\mathcal{N}(\mu_0, \tau^2)$ and the inverse Gamma distribution $IG(a, b)$, then

$$\left\{ \begin{array}{l} \mu | \sigma^2, \mathbf{x} \sim \mathcal{N} \left(\frac{\sigma^2}{\sigma^2 + n\tau^2} \mu_0 + \frac{n\tau^2}{\sigma^2 + n\tau^2} \bar{x}, \frac{\sigma^2 \tau^2}{\sigma^2 + n\tau^2} \right) \\ \sigma^2 | \mu, \mathbf{x} \sim IG \left(\frac{n}{2} + a, \frac{1}{2} \sum_{i=1}^n [x_i - \mu]^2 + b \right) \end{array} \right.$$

More generally, we need the conditional distribution of $\theta_k | \boldsymbol{\theta}_{-k}, \mathbf{x}$, for all k .

Gibbs Sampler



Gibbs Sampler

Example Consider some vector $\mathbf{X} = (X_1, \dots, X_d)$ with independent components, $X_i \sim \mathcal{E}(\lambda_i)$. To sample from \mathbf{X} given $\mathbf{X}^\top \mathbf{1} > s$ for some $s > 0$:

start with some starting point \mathbf{x}_0 such that $\mathbf{x}_0^\top \mathbf{1} > s$

pick up (randomly) $i \in \{1, \dots, d\}$

X_i given $X_i > s - \mathbf{x}_{(-i)}^\top \mathbf{1}$ has an Exponential distribution $\mathcal{E}(\lambda_i)$

draw $Y \sim \mathcal{E}(\lambda_i)$ and set $x_i = y + (s - \mathbf{x}_{(-i)}^\top \mathbf{1})_+$ until $\mathbf{x}_{(-i)}^\top \mathbf{1} + x_i > s$

JAGS and STAN

Martyn Plummer developed **JAGS** *Just another Gibbs sampler* in 2007 (stable since 2013). It is an open-source, enhanced, cross-platform version of an earlier engine BUGS (Bayesian inference Using Gibbs Sampling).

STAN is a newer tool that uses the Hamiltonian Monte Carlo (HMC) sampler.

HMC uses information about the derivative of the posterior probability density to improve the algorithm. These derivatives are supplied by algorithm differentiation in C/C++ codes.

MCMC and Claims Reserving

Consider the following (cumulated) triangle, $\{C_{i,j}\}$,

	0	1	2	3	4	5
0	3209	4372	4411	4428	4435	4456
1	3367	4659	4696	4720	4730	4752.4
2	3871	5345	5398	5420	5430.1	5455.8
3	4239	5917	6020	6046.1	6057.4	6086.1
4	4929	6794	6871.7	6901.5	6914.3	6947.1
5	5217	7204.3	7286.7	7318.3	7331.9	7366.7

λ_j	1.3809	1.0114	1.0043	1.0018	1.0047
σ_j	0.7248	0.3203	0.04587	0.02570	0.02570

A Bayesian version of Chain Ladder

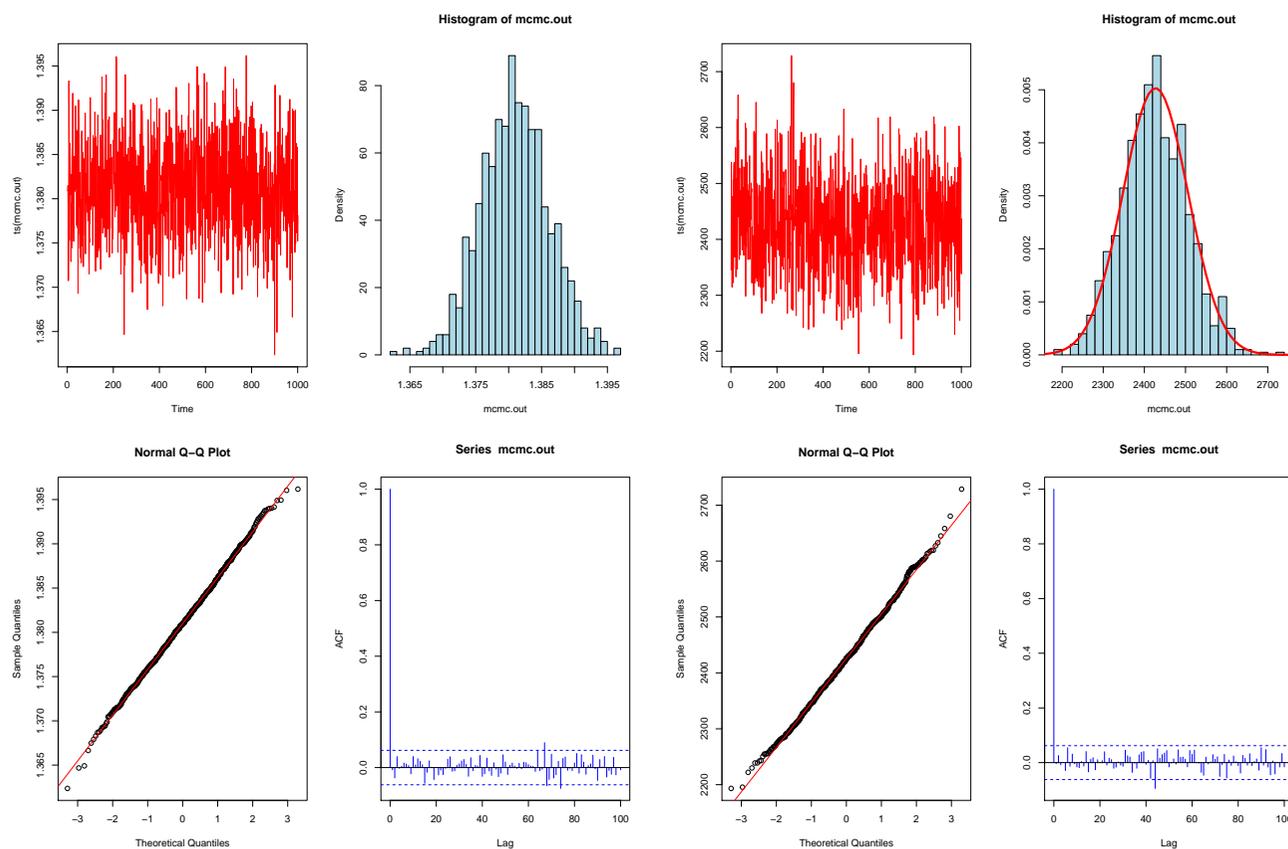
	0	1	2	3	4	5
0		1.362418	1.008920	1.003854	1.001581	1.004735
1		1.383724	1.007942	1.005111	1.002119	
2		1.380780	1.009916	1.004076		
3		1.395848	1.017407			
4		1.378373				

λ_j	1.3809	1.0114	1.0043	1.0018	1.0047
σ_j	0.7248	0.3203	0.04587	0.02570	0.02570

Assume that $\lambda_{i,j} \sim \mathcal{N}\left(\mu_j, \frac{\tau_j}{C_{i,j}}\right)$.

We can use Gibbs sampler to get the distribution of the transition factors, as well as a distribution for the reserves,

A Bayesian version of Chain Ladder



A Bayesian analysis of the Poisson Regression Model

In a Poisson regression model, we have a sample

$$(\mathbf{x}, \mathbf{y}) = \{(x_i, y_i)\},$$

$$y_i \sim \mathcal{P}(\mu_i) \text{ with } \log \mu_i = \beta_0 + \beta_1 x_i.$$

In the Bayesian framework, β_0 and β_1 are random variables.

Other alternatives to classical statistics

Consider a regression problem, $\mu(x) = \mathbb{E}(Y|X = x)$, and assume that smoothed splines are used,

$$\mu(x) = \sum_{i=1}^k \beta_j h_j(x)$$

Let \mathbf{H} be the $n \times k$ matrix, $\mathbf{H} = [h_j(x_i)] = [\mathbf{h}(x_i)]$, then $\hat{\boldsymbol{\beta}} = (\mathbf{H}^\top \mathbf{H})^{-1} \mathbf{H}^\top \mathbf{y}$, and

$$\widehat{\text{se}}(\hat{\mu}(x)) = [\mathbf{h}(x)^\top (\mathbf{H}^\top \mathbf{H})^{-1} \mathbf{h}(x)]^{\frac{1}{2}} \hat{\sigma}$$

With a Gaussian assumption on the residuals, we can derive (approximated) confidence bands for predictions $\hat{\mu}(x)$.

Bayesian interpretation of the regression problem

Assume here that $\boldsymbol{\beta} \sim \mathcal{N}(\mathbf{0}, \tau \boldsymbol{\Sigma})$ as the priori distribution for $\boldsymbol{\beta}$.

Then, if $(\mathbf{x}, \mathbf{y}) = \{(x_i, y_i), i = 1, \dots, n\}$, the posterior distribution of $\mu(x)$ will be Gaussian, with

$$\mathbb{E}(\mu(x) | \mathbf{x}, \mathbf{y}) = \mathbf{h}(x)^\top \left(\mathbf{H}^\top \mathbf{H} + \frac{\sigma^2}{\tau} \boldsymbol{\Sigma}^{-1} \right)^{-1} \mathbf{H}^\top \mathbf{y}$$

$$\text{cov}(\mu(x), \mu(x') | \mathbf{x}, \mathbf{y})$$

$$= \mathbf{h}(x)^\top \left(\mathbf{H}^\top \mathbf{H} + \frac{\sigma^2}{\tau} \boldsymbol{\Sigma}^{-1} \right)^{-1} \mathbf{h}(x') \sigma^2$$

Example $\boldsymbol{\Sigma} = \mathbb{I}$

Bootstrap strategy

Assume that $Y = \mu(x) + \varepsilon$, and based on the estimated model, generate pseudo observations, $y_i^* = \hat{\mu}(x_i) + \hat{\varepsilon}_i^*$.

Based on $(\mathbf{x}, \mathbf{y}^*) = \{(x_i, y_i^*), i = 1, \dots, n\}$, derive the estimator $\hat{\mu}^*(\cdot)$

(and repeat)

Observe that the bootstrap is the Bayesian case, when $\tau \rightarrow \infty$.

Part 5. Data, Models & Actuarial Science (some sort of conclusion)



The Privacy-Utility Trade-Off

In Massachusetts, the Group Insurance Commission (GIC) is responsible for purchasing health insurance for state employees

GIC has to publish the data: **GIC**(zip, date of birth, sex, diagnosis, procedure, ...)

Sweeney paid \$20 and bought the voter registration list for Cambridge Massachusetts, **VOTER**(name, party, ..., zip, date of birth, sex)

William Weld (former governor) lives in Cambridge, hence is in **VOTER**

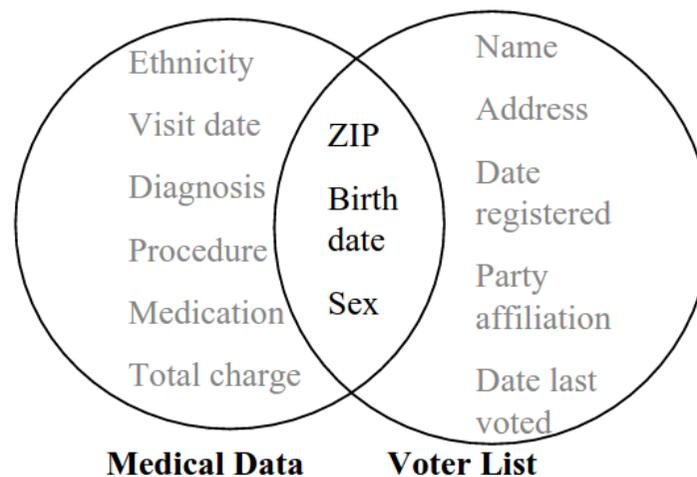


Figure 1 Linking to re-identify data

The Privacy-Utility Trade-Off

- 6 people in **VOTER** share his date of birth
- only 3 of them were man (same sex)
- Weld was the only one in that zip
- Sweeney learned Weld's medical records



All systems worked as specified, yet an important data was leaked.

“87% of Americans are uniquely identified by their zip code, gender and birth date”, see [Sweeney \(2000\)](#).

A dataset is considered k -anonymous if the information for each person contained in the release cannot be distinguished from at least $k - 1$ individuals whose information also appear in the release

manière des caractéristiques Ω de l'assuré, et lui réclame donc une prime pure de montant $\mathbb{E}[S]$, la même que celle qu'il réclame à tous les assurés du portefeuille. Dans ce cas, la situation est telle que présentée au Tableau 3.7.

	Assurés	Assureur
Dépense	$\mathbb{E}[S]$	$S - \mathbb{E}[S]$
Dépense moyenne	$\mathbb{E}[S]$	0
Variance	0	$\mathbb{V}[S]$

TAB. 3.7 – Situation des assurés et de l'assureur en l'absence de segmentation.

L'assureur prend donc l'entièreté de la variance des sinistres $\mathbb{V}[S]$ à sa charge, que celle-ci soit due à l'hétérogénéité du portefeuille, ou à la variabilité intrinsèque des montants des sinistres.

Transfert de risque en information complète

À l'autre extrême, supposons que l'assureur incorpore toute l'information Ω dans la tarification. On serait alors dans la situation décrite au Tableau 3.8.

	Assurés	Assureur
Dépense	$\mathbb{E}[S \Omega]$	$S - \mathbb{E}[S \Omega]$
Dépense moyenne	$\mathbb{E}[S]$	0
Variance	$\mathbb{V}[\mathbb{E}[S \Omega]]$	$\mathbb{V}[S - \mathbb{E}[S \Omega]]$

TAB. 3.8 – Situation des assurés et de l'assureur dans le cas où la segmentation est opérée sur base de Ω .

Contrairement au cas précédent, la prime payée par un assuré prélevé au hasard dans le portefeuille est à présent une variable aléatoire: $\mathbb{E}[S|\Omega]$ dépend des caractéristiques Ω de cet assuré. Comme la variable aléatoire $S - \mathbb{E}[S|\Omega]$ est centrée, le risque assumé par l'assureur la variance du résultat financier de l'opération d'assurance, i.e.

$$\mathbb{V}[S - \mathbb{E}[S|\Omega]] = \mathbb{E}[(S - \mathbb{E}[S|\Omega])^2]$$

No segmentation

	Insured	Insurer
Loss	$\mathbb{E}[S]$	$S - \mathbb{E}[S]$
Average Loss	$\mathbb{E}[S]$	0
Variance	0	$\text{Var}[S]$

Perfect Information: Ω observable

	Insured	Insurer
Loss	$\mathbb{E}[S \Omega]$	$S - \mathbb{E}[S \Omega]$
Average Loss	$\mathbb{E}[S]$	0
Variance	$\text{Var}[\mathbb{E}[S \Omega]]$	$\text{Var}[S - \mathbb{E}[S \Omega]]$

$$\text{Var}[S] = \underbrace{\mathbb{E}[\text{Var}[S|\Omega]]}_{\rightarrow \text{insurer}} + \underbrace{\text{Var}[\mathbb{E}[S|\Omega]]}_{\rightarrow \text{insured}}.$$

3.8. La prime pure en univers segmenté

177

On assiste dans ce cas à un partage de la variance totale de S (c'est-à-dire du risque) entre les assurés et l'assureur, matérialisé par la formule

$$\mathbb{V}[S] = \underbrace{\mathbb{E}[\mathbb{V}[S|\Omega]]}_{\rightarrow \text{assureur}} + \underbrace{\mathbb{V}[\mathbb{E}[S|\Omega]]}_{\rightarrow \text{assurés}}.$$

Ainsi, lorsque toutes les variables pertinentes Ω ont été prises en compte, l'intervention de l'assureur se limite à la part des sinistres due exclusivement au hasard; en effet, $\mathbb{V}[S|\Omega]$ représente les fluctuations de S dues au seul hasard. Dans cette situation idéale, l'assureur mutualise le risque et il n'y a donc aucune solidarité induite entre les assurés du portefeuille: chacun paie en fonction de son propre risque.

Transfert des risques en information partielle

Bien entendu, la situation décrite au paragraphe précédent est purement théorique puisque parmi les variables explicatives Ω nombreuses sont celles qui ne peuvent pas être observées par l'assureur. En assurance automobile par exemple, l'assureur ne peut pas observer la vitesse à laquelle roule l'assuré, son agressivité au volant, ni le nombre de kilomètres qu'il parcourt chaque année². Dès lors, l'assureur ne peut utiliser qu'un sous-ensemble \mathbf{X} des variables explicatives contenues dans Ω , i.e. $\mathbf{X} \subset \Omega$. La situation est alors semblable à celle décrite au Tableau 3.9.

	Assuré	Assureur
Dépense	$\mathbb{E}[S \mathbf{X}]$	$S - \mathbb{E}[S \mathbf{X}]$
Dépense moyenne	$\mathbb{E}[S]$	0
Variance	$\mathbb{V}[\mathbb{E}[S \mathbf{X}]]$	$\mathbb{E}[\mathbb{V}[S \mathbf{X}]]$

TAB. 3.9 – Situation de l'assuré et de l'assureur dans le cas où la segmentation est opérée sur base de $\mathbf{X} \subset \Omega$.

Il est intéressant de constater que

$$\begin{aligned} \mathbb{E}[\mathbb{V}[S|\mathbf{X}]] &= \mathbb{E}[\mathbb{E}[\mathbb{V}[S|\Omega]|\mathbf{X}]] + \mathbb{E}[\mathbb{V}[\mathbb{E}[S|\Omega]|\mathbf{X}]] \\ &= \underbrace{\mathbb{E}[\mathbb{V}[S|\Omega]]}_{\text{mutualisation}} + \underbrace{\mathbb{E}[\mathbb{V}[\mathbb{E}[S|\Omega]|\mathbf{X}]]}_{\text{solidarité}}. \end{aligned} \quad (3.22)$$

Non-Perfect Information: $\mathbf{X} \subset \Omega$ is observable

	Insured	Insurer
Loss	$\mathbb{E}[S \mathbf{X}]$	$S - \mathbb{E}[S \mathbf{X}]$
Average Loss	$\mathbb{E}[S]$	0
Variance	$\text{Var}[\mathbb{E}[S \mathbf{X}]]$	$\mathbb{E}[\text{Var}[S \mathbf{X}]]$

$$\begin{aligned} \mathbb{E}[\text{Var}[S|\mathbf{X}]] &= \mathbb{E}[\mathbb{E}[\text{Var}[S|\Omega]|\mathbf{X}]] \\ &+ \mathbb{E}[\text{Var}[\mathbb{E}[S|\Omega]|\mathbf{X}]] \\ &= \underbrace{\mathbb{E}[\text{Var}[S|\Omega]]}_{\text{pooling}} \\ &+ \underbrace{\mathbb{E}\{\text{Var}[\mathbb{E}[S|\Omega]|\mathbf{X}]\}}_{\text{solidarity}}. \end{aligned}$$

SEGMENTATION ET MUTUALISATION LES DEUX FACES D'UNE MÊME PIÈCE ?

Arthur Charpentier

Professeur à l'Université du Québec, Montréal

Michel Denuit

Professeur à l'Université catholique de Louvain

Romuald Elie

Professeur à l'Université de Marne-la-Vallée

L'assurance repose fondamentalement sur l'idée que la mutualisation des risques entre des assurés est possible. Cette mutualisation, qui peut être vue comme une relecture actuarielle de la loi des grands nombres, n'a de sens qu'au sein d'une population de risques « homogènes » [Charpentier, 2011]. Cette condition (actuarielle) impose aux assureurs de segmenter, ce que confirment plusieurs travaux économiques ⁽¹⁾. Avec l'explosion du nombre de données, et donc de variables tarifaires possibles, certains assureurs évoquent l'idée d'un tarif individuel, semblant remettre en cause l'idée même de mutualisation des risques. Entre cette force qui pousse à segmenter et la force de rappel qui tend (pour des raisons sociales mais aussi actuarielles, ou au moins de robustesse statistique ⁽²⁾) à imposer une solidarité minimale entre les assurés, quel équilibre va en résulter dans un contexte de forte concurrence entre les sociétés d'assurance ?

Tarifcation sans segmentation

Sans segmentation, le « prix juste » d'un risque est l'espérance mathématique de la charge annuelle. C'est l'idée du théorème fondamental de la valorisation actuarielle : en moyenne, la somme des primes doit permettre d'indemniser l'intégralité des sinistres survenus dans

l'année. Afin d'illustrer les différents aspects de la construction du tarif et ses conséquences, on va utiliser les données présentées dans le tableau 1 (voir p. xx), qui indique la fréquence annuelle de sinistres.

Les facteurs de risque sont ici le lieu d'habitation et l'âge de l'assuré, et on observe la fréquence de sinistre par classe. Le coût unitaire, supposé fixe, équivaut à 1 000 euros. La prime pure est alors $E[S] = 1\,000 \times E[N]$. Dans cet exemple, la prime pure sans segmentation sera de 82,30 euros.

Simple model $\Omega = \{X_1, X_2\}$.

Four Models

$$\left\{ \begin{array}{l} \hat{m}_0(\mathbf{x}_1, \mathbf{x}_2) = \mathbb{E}[S] \\ \hat{m}_1(\mathbf{x}_1, \mathbf{x}_2) = \mathbb{E}[S | X_1 = \mathbf{x}_1] \\ \hat{m}_2(\mathbf{x}_1, \mathbf{x}_2) = \mathbb{E}[S | X_2 = \mathbf{x}_2] \\ \hat{m}_{12}(\mathbf{x}_1, \mathbf{x}_2) = \mathbb{E}[S | X_1 = \mathbf{x}_1, X_2 = \mathbf{x}_2] \end{array} \right.$$

Market Competition

Decision Rule: the insured selects the **cheapeast premium**,

	A	B	C	D	E	F
	787.93	706.97	1032.62	907.64	822.58	603.83
	170.04	197.81	285.99	212.71	177.87	265.13
	473.15	447.58	343.64	410.76	414.23	425.23
	337.98	336.20	468.45	339.33	383.55	672.91

Market Competition

Decision Rule: the insured selects randomly from the **three cheapest premium**

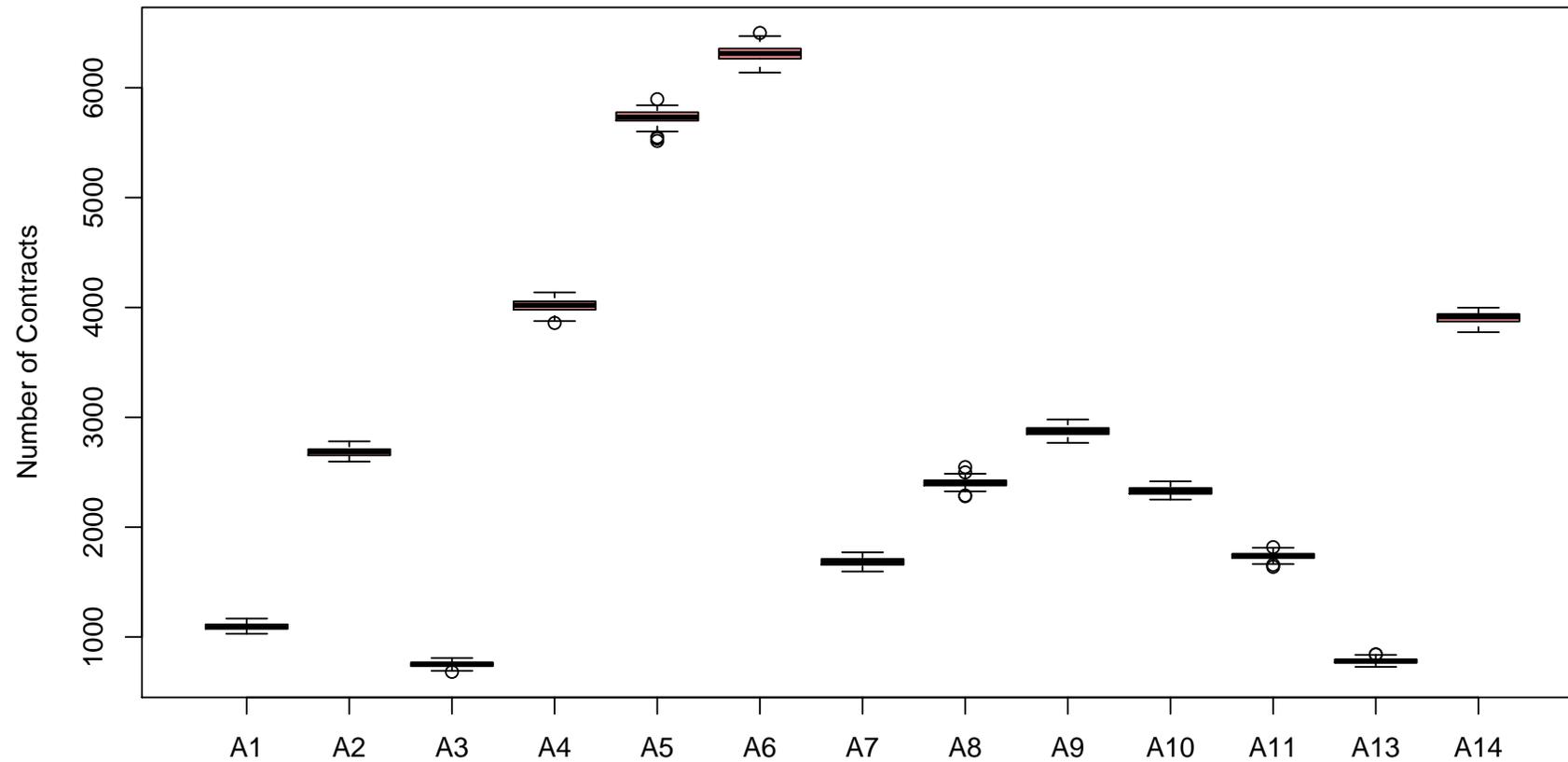
	A	B	C	D	E	F
	787.93	706.97	1032.62	907.64	822.58	603.83
	170.04	197.81	285.99	212.71	177.87	265.13
	473.15	447.58	343.64	410.76	414.23	425.23
	337.98	336.20	468.45	339.33	383.55	672.91

Market Competition

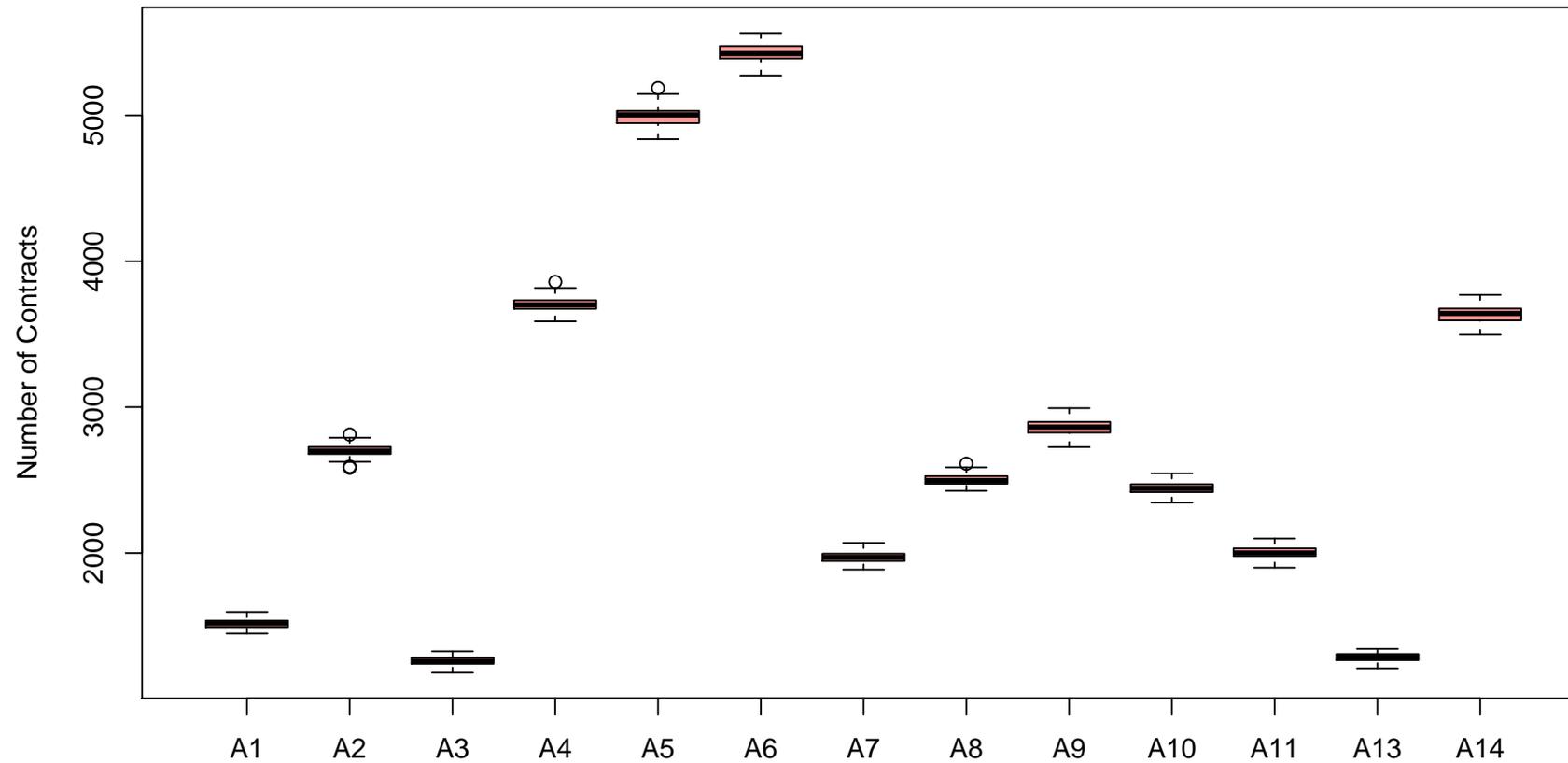
Decision Rule: the insured were assigned randomly to some insurance company for year $n - 1$. For year n , they stay with their company if the premium is one of the three cheapest premium, if not, random choice among the four

	A	B	C	D	E	F
	787.93	706.97	1032.62	907.64	822.58	603.83
	170.04	197.81	285.99	212.71	177.87	265.13
	473.15	447.58	343.64	410.76	414.23	425.23
	337.98	336.20	468.45	339.33	383.55	672.91

Market Shares (rule 2)

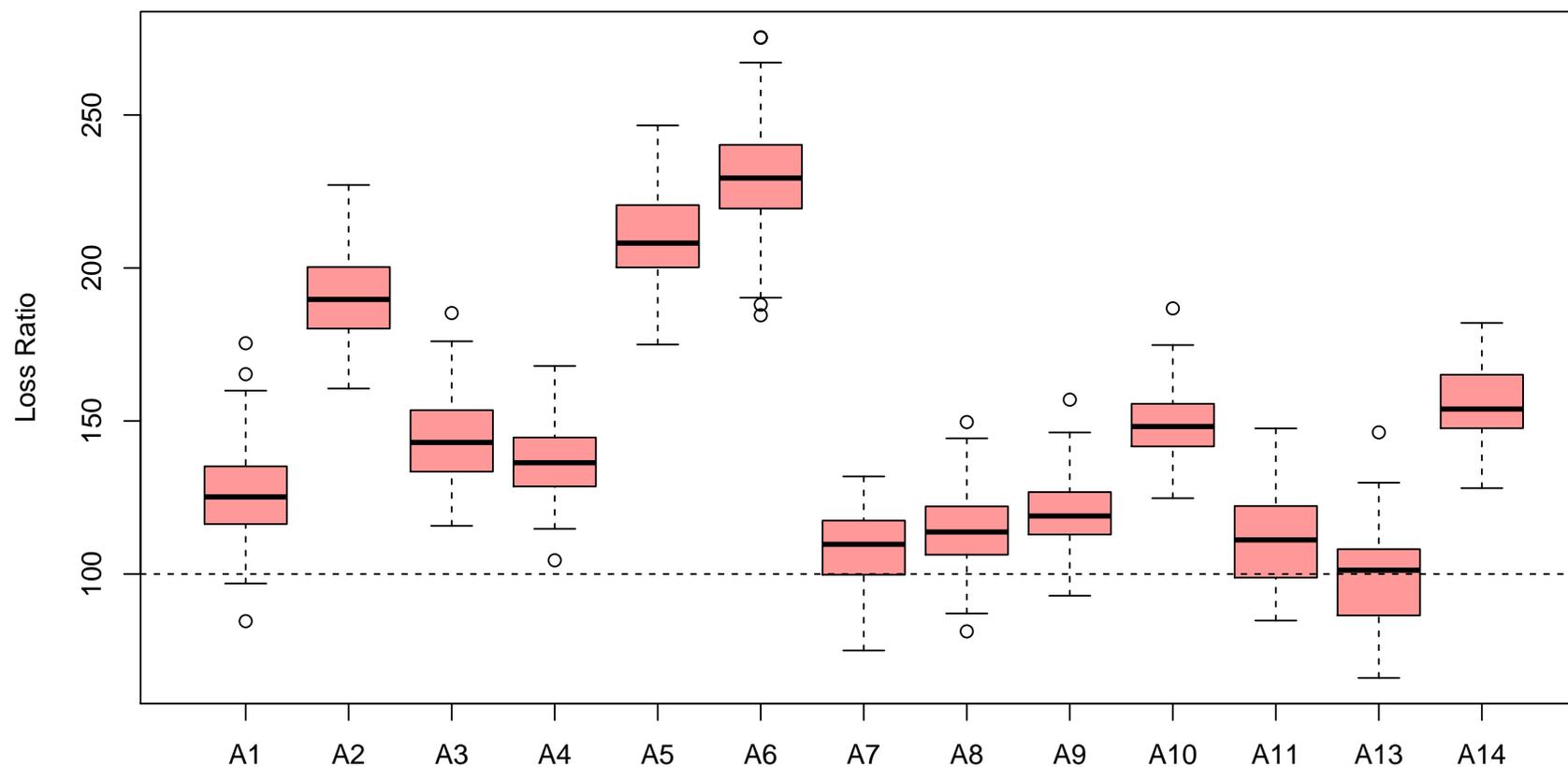


Market Shares (rule 3)



Loss Ratio, Loss / Premium (rule 2)

Market Loss Ratio $\sim 154\%$.

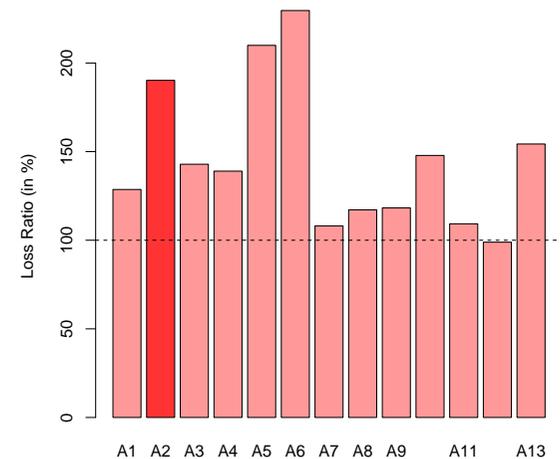
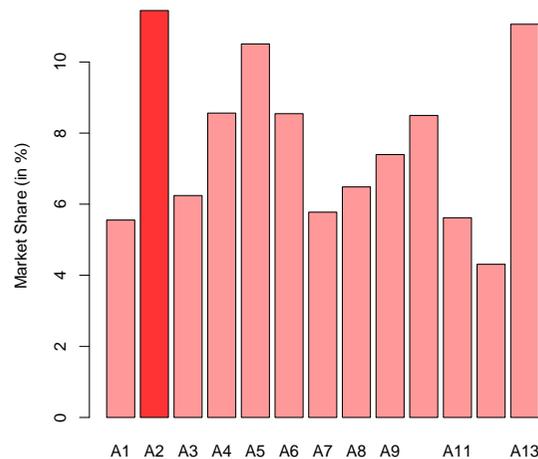
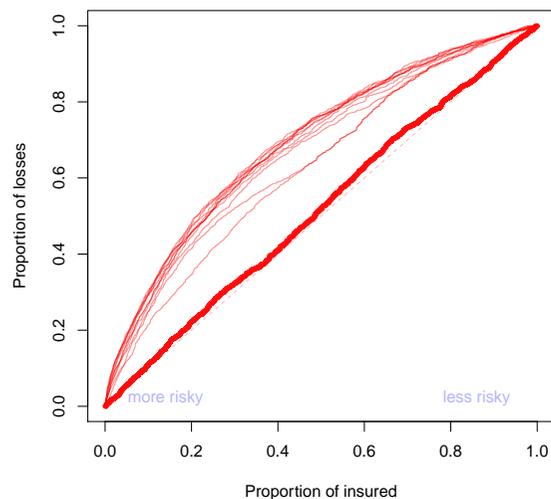


Insurer A2

No segmentation, unique premium

Remark on normalized premiums,

$$\pi_2 = m_2(\mathbf{x}_i) = \frac{1}{n} \sum_{i=1}^n m_j(\mathbf{x}_i) \quad \forall j$$



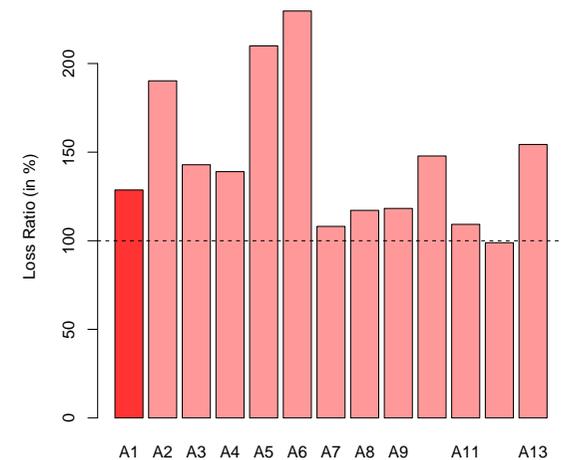
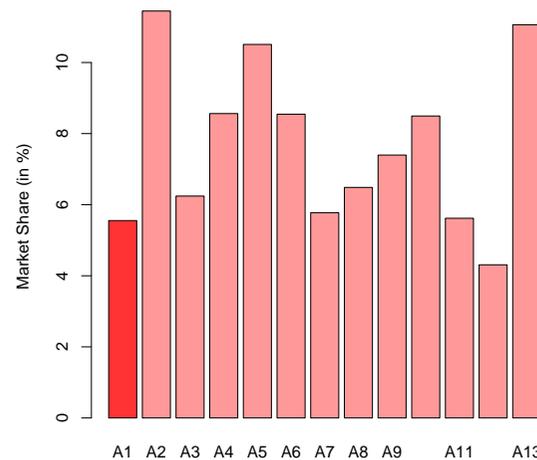
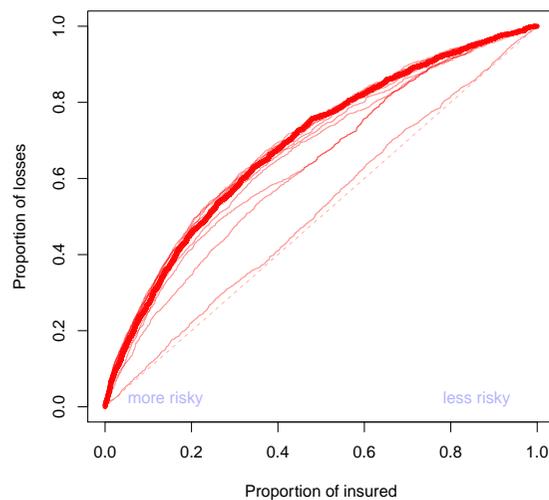
Insured A1

GLM, frequency material / bodily injury, individual losses material

Ages in classes [18-30], [30-45], [45-60] and [60+], crossed with occupation

Manual smoothing, SAS and Excel

Actuaries in a Mutual Fund (in France)



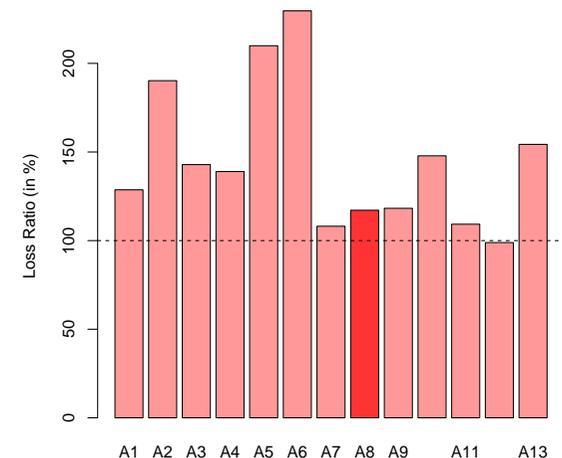
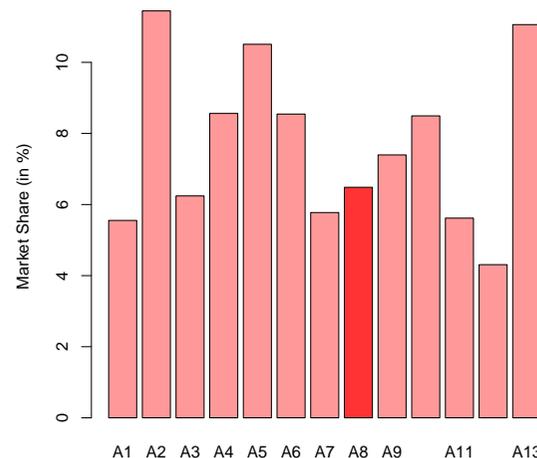
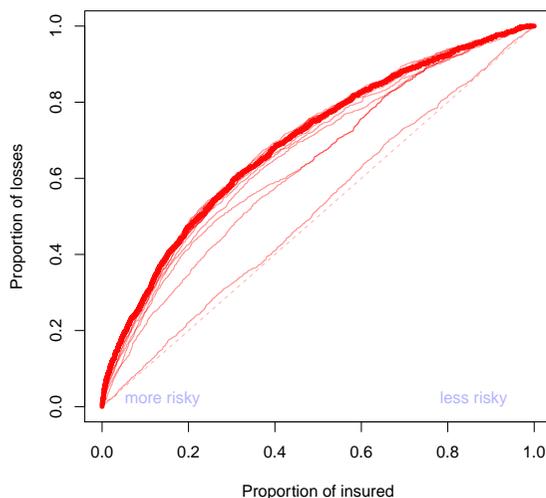
Insurer A8/A9

GLM, frequency and losses, without major losses ($>15k$)

Age-gender interaction

Use of a commercial pricing software

Actuary in a French Mutual Fund

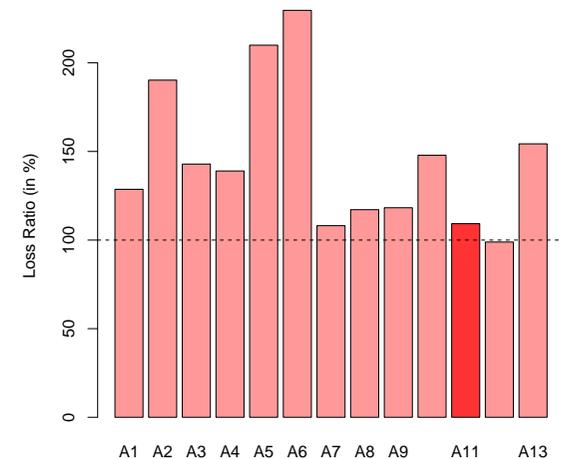
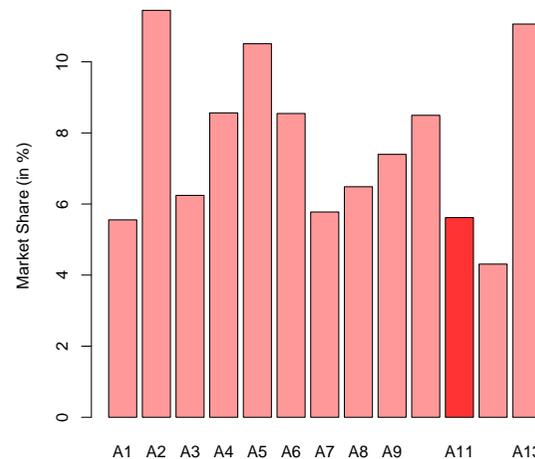
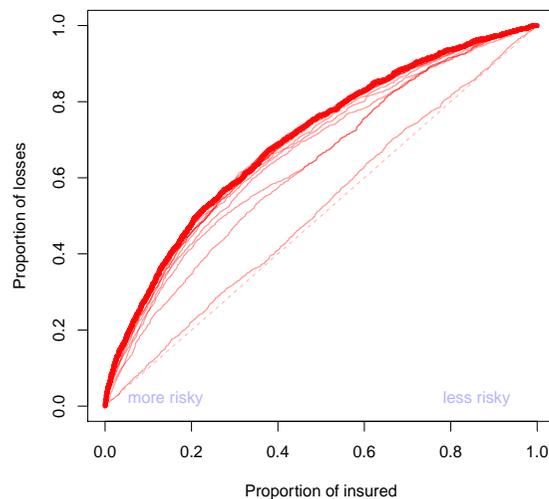


Insurer A11

All features, but one XGBoost (gradient boosting)

Correction for negative premiums

Coded in Python actuary in an insurance company.

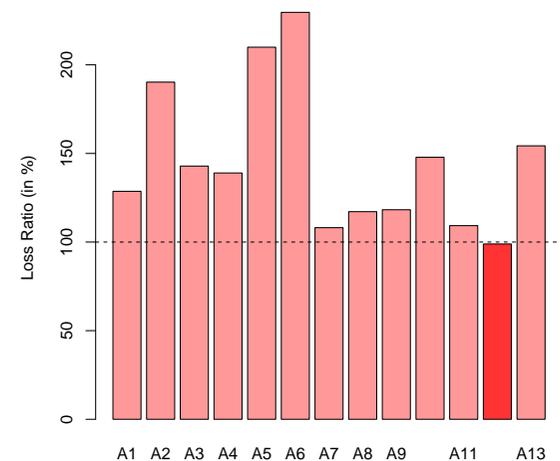
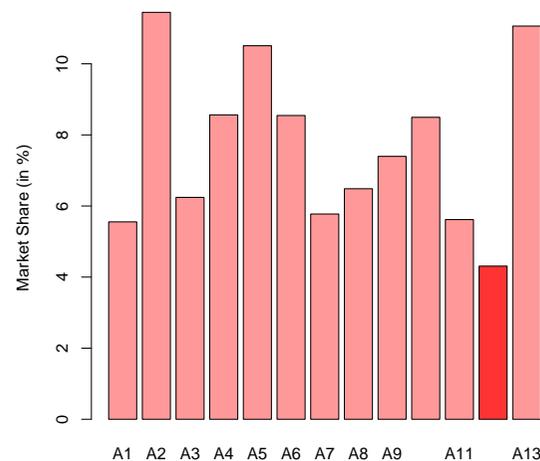
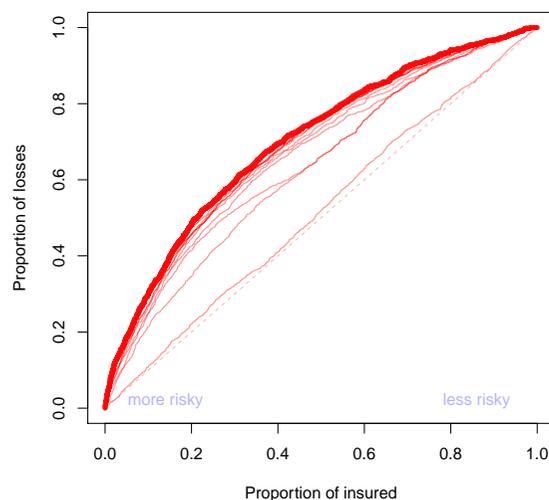


Insurer A12

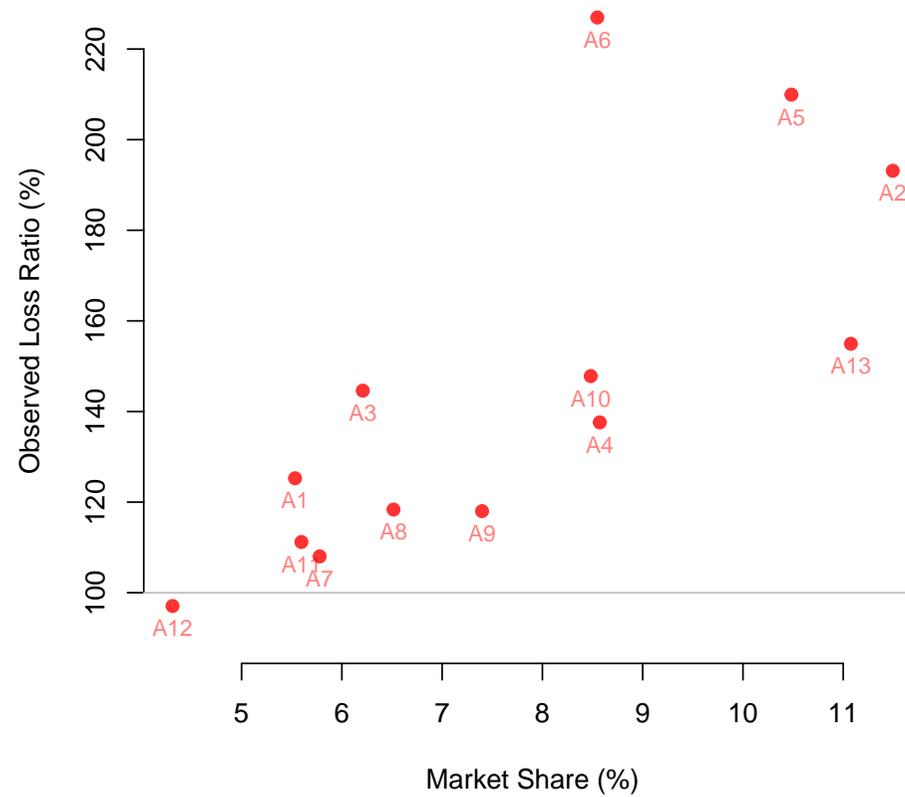
All features, use of two XGBoost (gradient boosting) models

Correction for negative premiums

Coded in R by an actuary in an Insurance company.



Back on the Pricing Game



Take-Away Conclusion

“People rarely succeed unless they have fun in what they are doing ” D. Carnegie

- on very small datasets, it is possible to use **Bayesian technique** to derive robust predictions,
- on extremely large datasets, it is possible to use ideas developed in **machine learning**, on regression models (e.g. bootstrapping and aggregating)
- all those techniques require **computational skills**

“the numbers have no way of speaking for themselves. We speak for them. ... Before we demand more of our data, we need to demand more of ourselves ” N. Silver, in **Silver (2012)**.

