

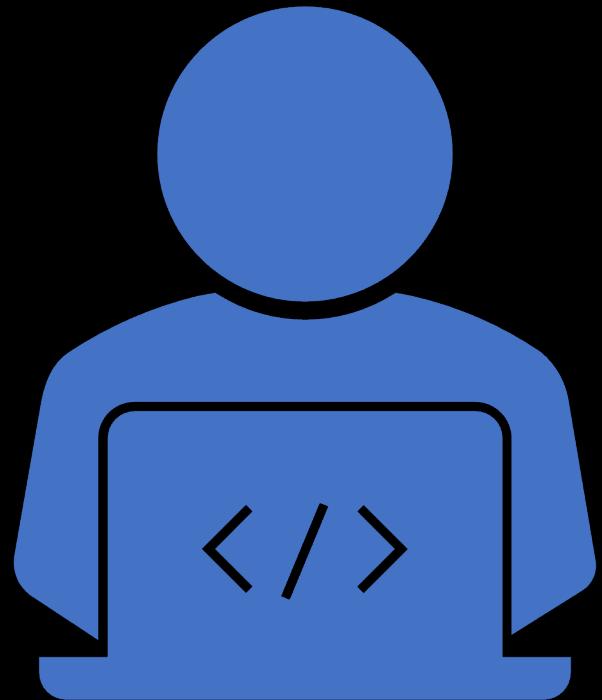
M1a: In the beginning

DSE 10200: Introduction to Data Science

Instructor: Michael Grossberg

What's up with “Data Science”?

Isn't it just statistics
computer science
ai
science

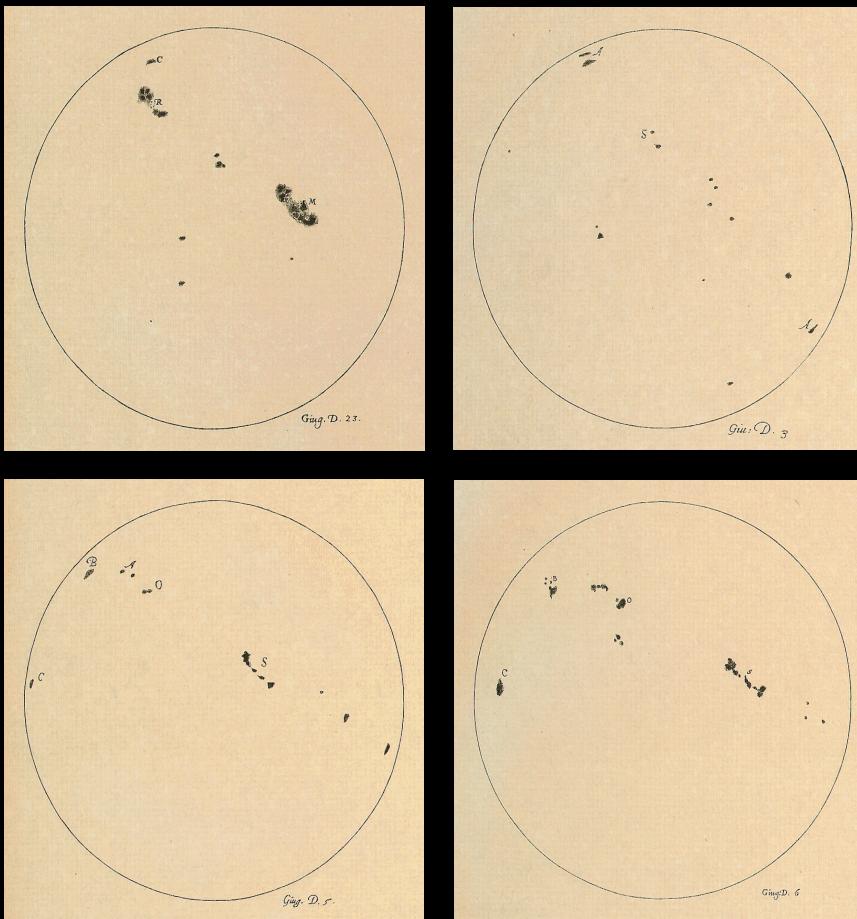


Throughout
Most of
Human
History



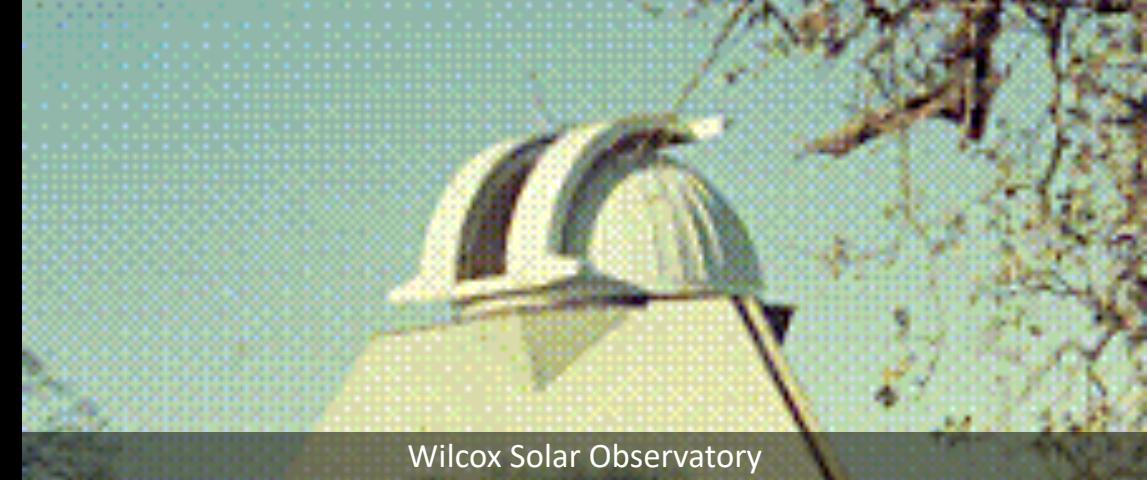
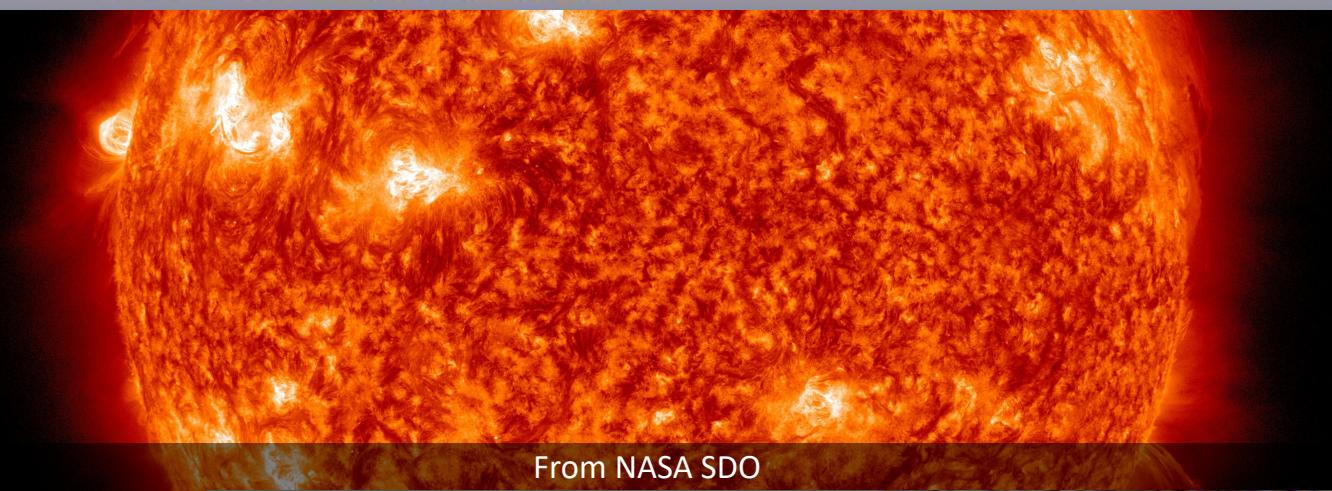
Data difficult to collect and store

Galileo's Hand-drawn Sunspots

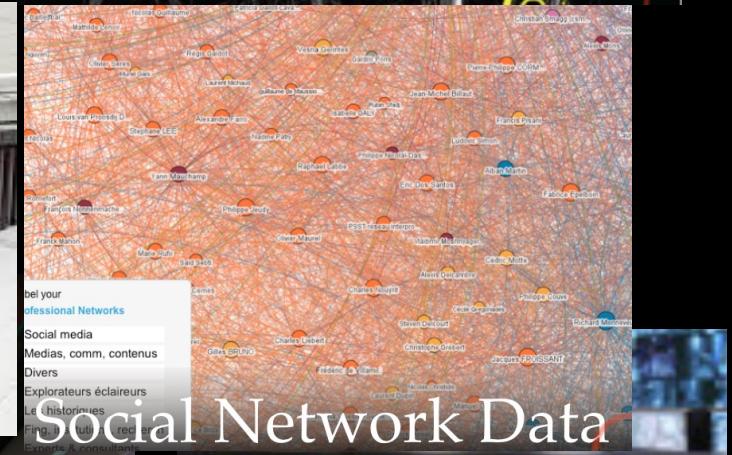
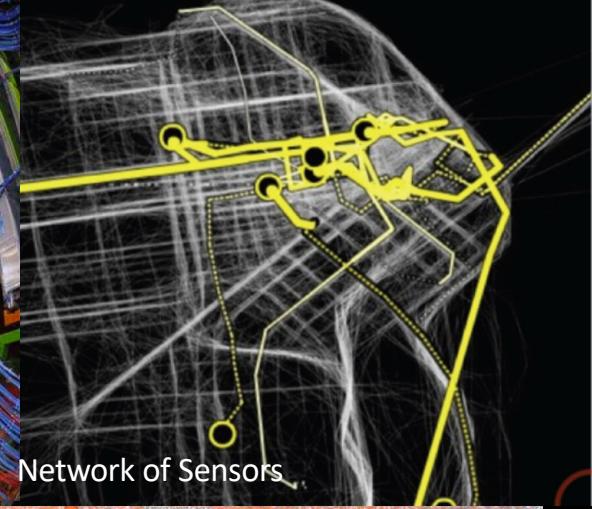
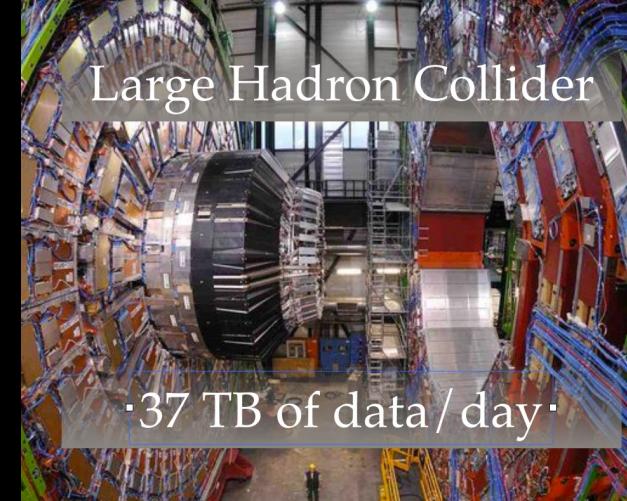


VS.

UltraHD video of
the Sun 24/7



Flood of Data



Every two days now we create
as much information as we did
from the dawn of civilization
up until 2003

- Eric Schmidt, former CEO Google, 2010



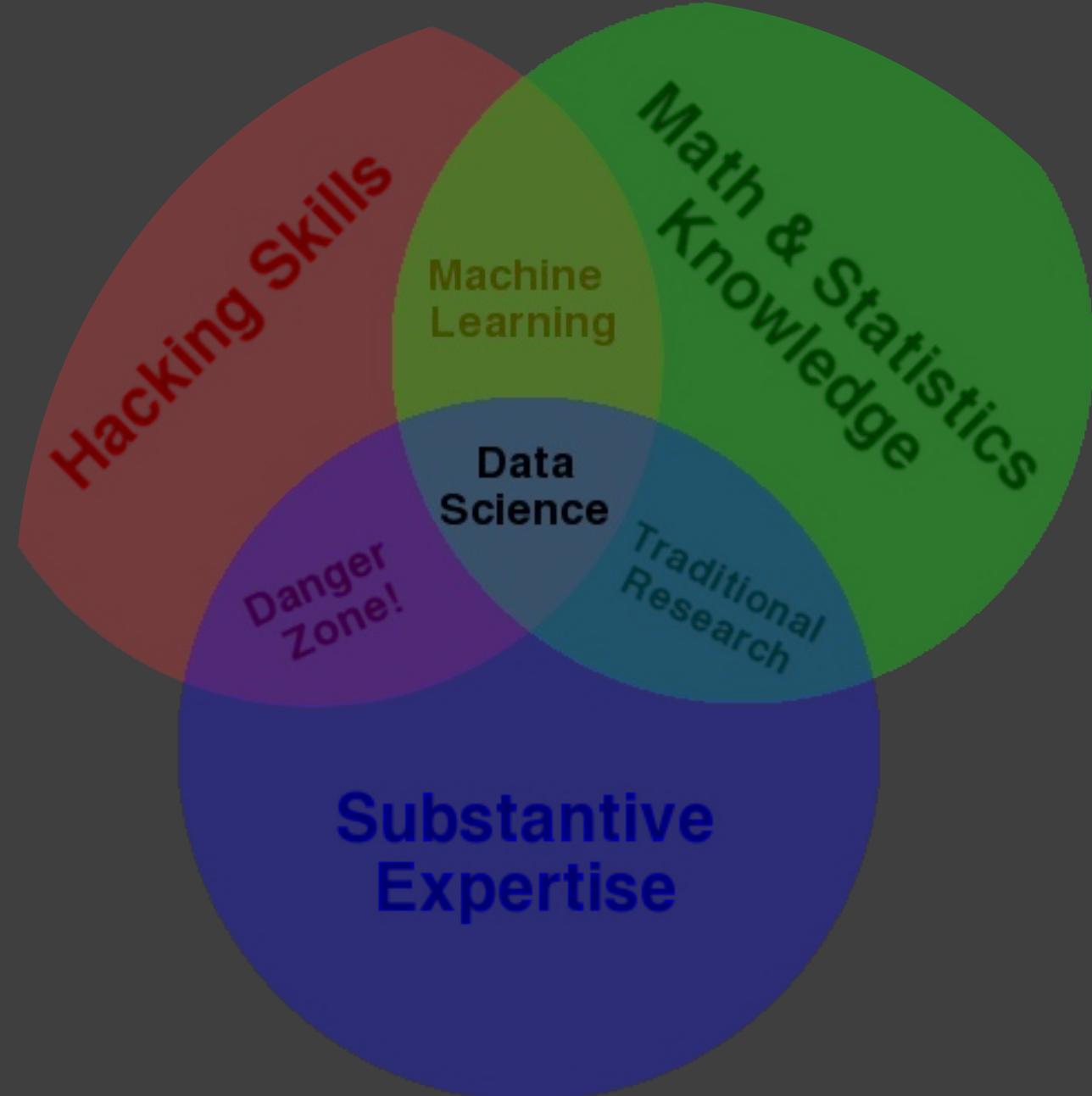
Data +
Computation



Is Data Science a new thing?

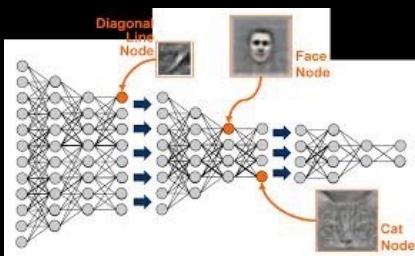
If things can combine to make something new

What is data science?



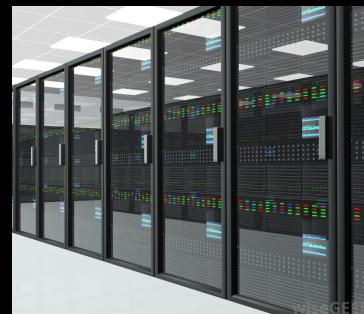
Coping with Data

Algorithms



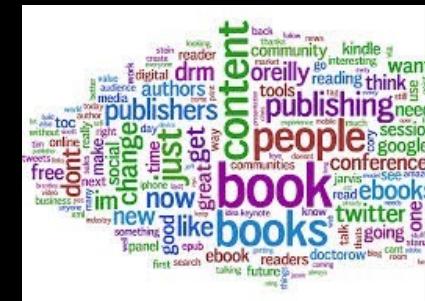
+

Computation



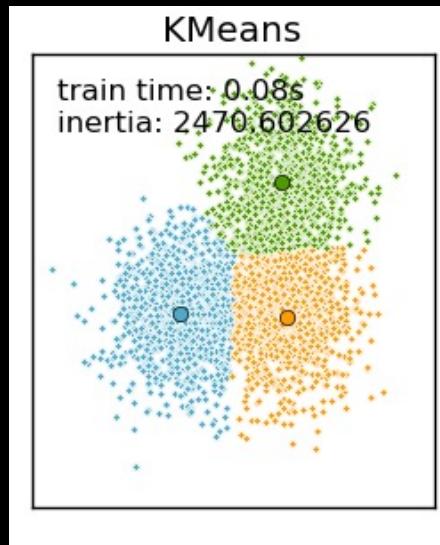
+

Visualization

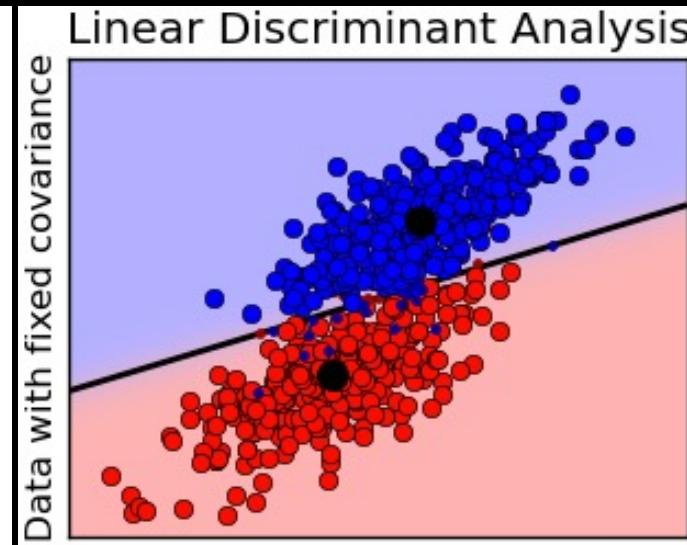


Core Subject Mater

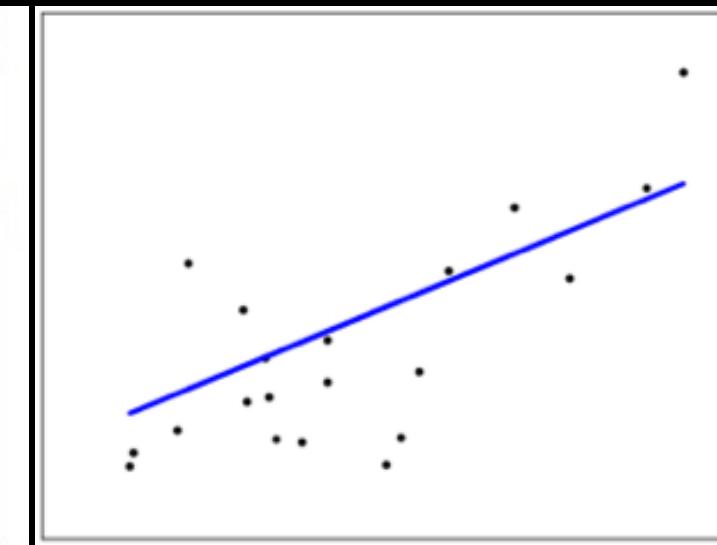
Machine Learning



Clustering Problem



Classification Problem



Regression Problem

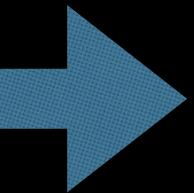
Goal: Explore structure of data and make predictions

Machine Learning in Engineering

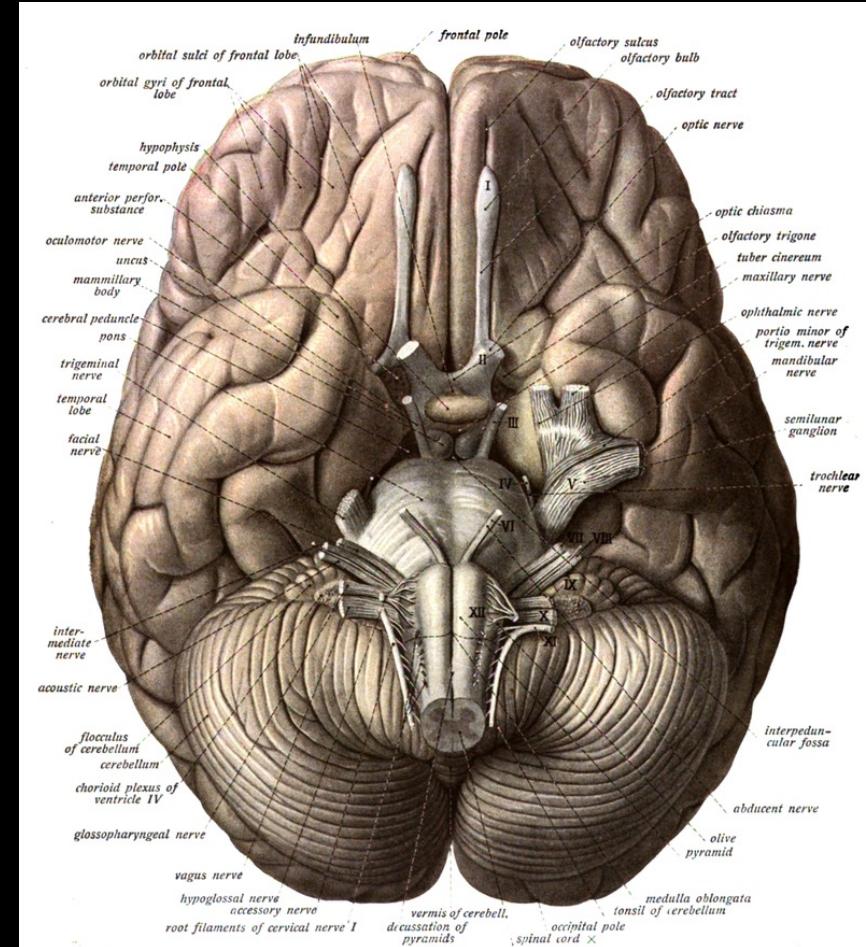
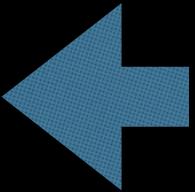


Data Visualization

Information



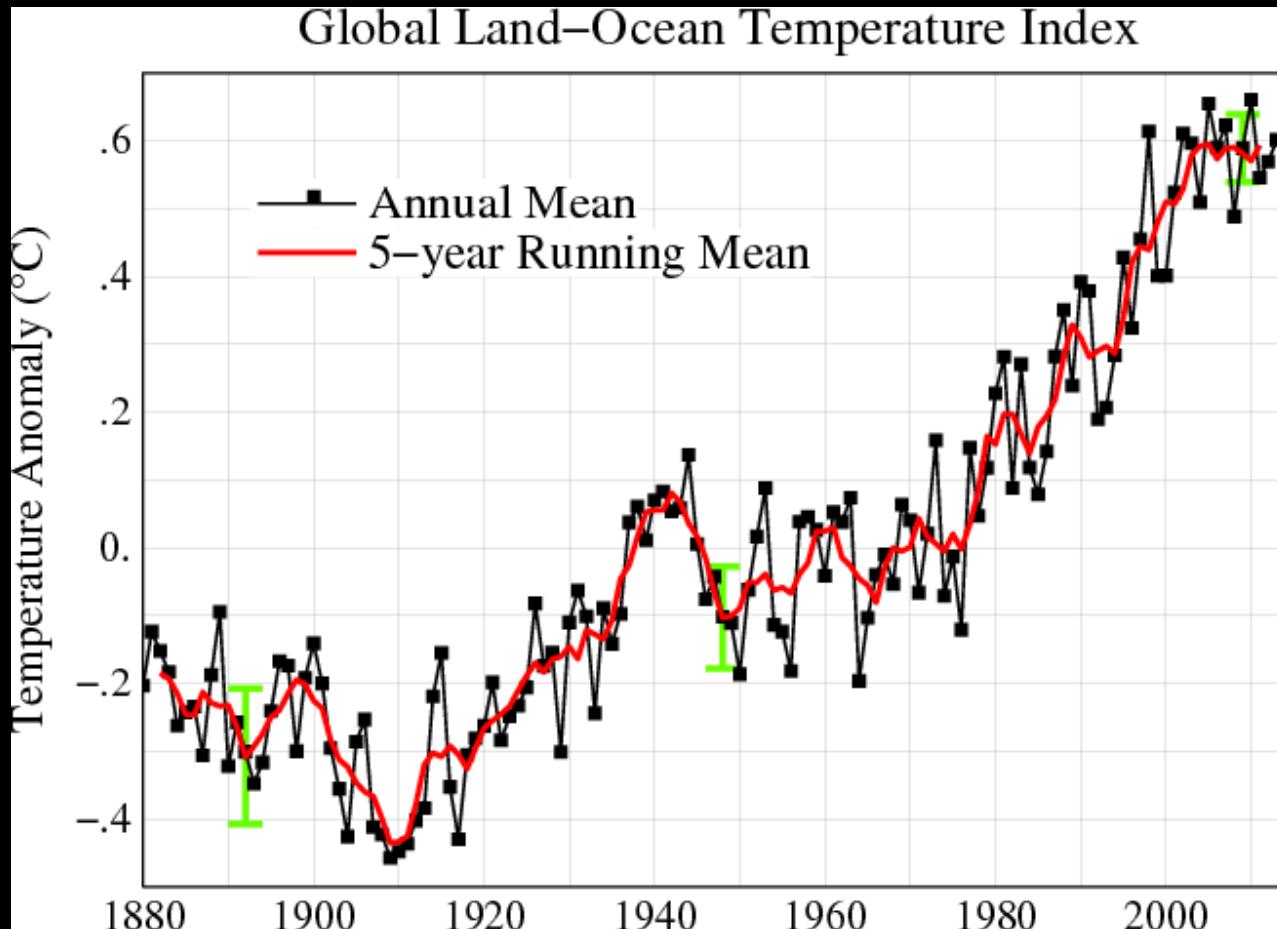
Understanding



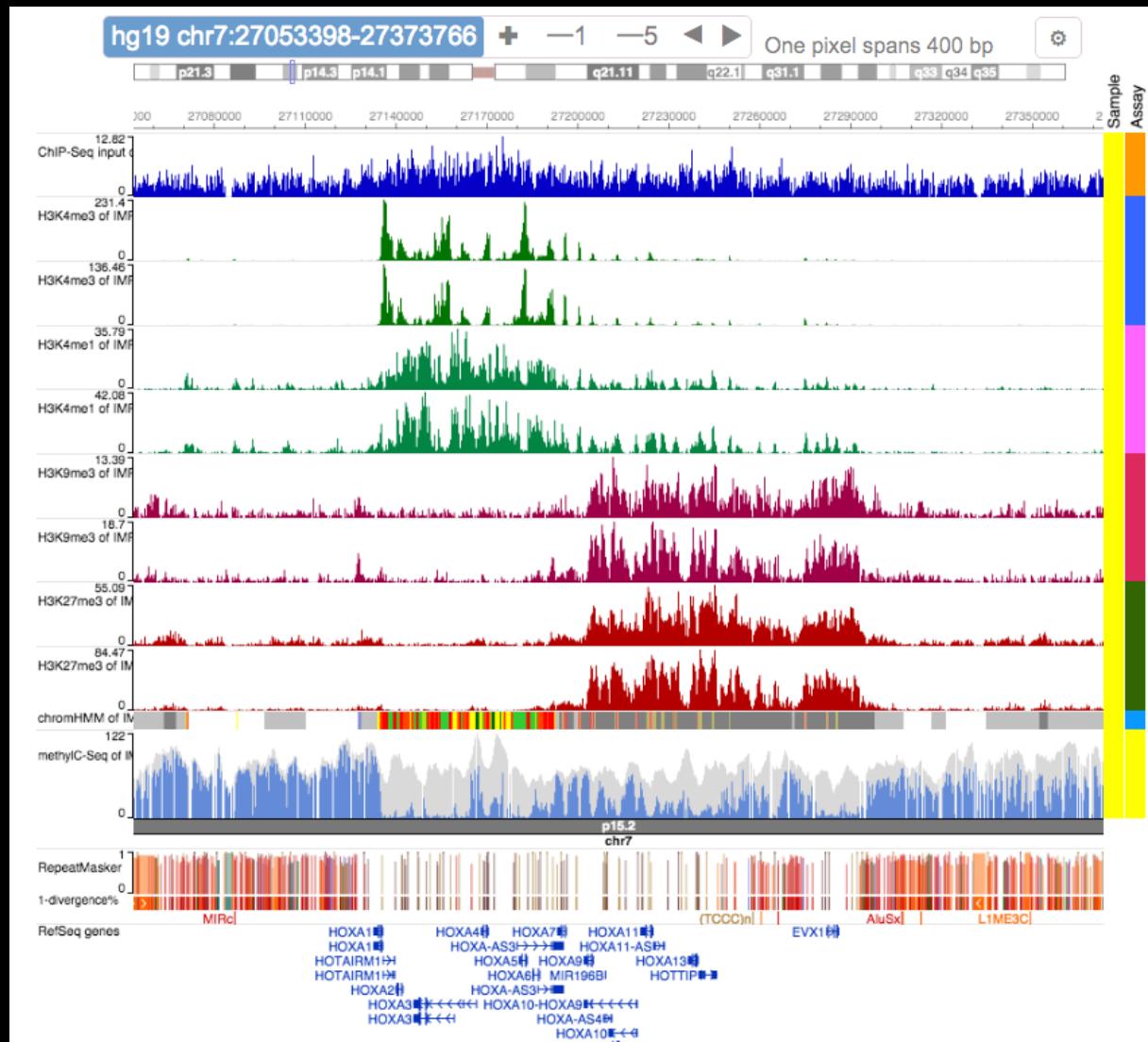
Data Visualization Global Temp

Year	Annual_Mean	5-year_Mean	1909	-0.46	-0.44	1942	0.05	0.08	1975	-0.01	0.02	2008	0.49	0.59
1880	-0.20	*	1910	-0.45	-0.43	1943	0.06	0.07	1976	-0.12	-0.00	2009	0.59	0.58
1881	-0.12	*	1911	-0.44	-0.43	1944	0.14	0.04	1977	0.15	0.04	2010	0.66	0.57
1882	-0.15	-0.19	1912	-0.40	-0.38	1945	0.01	0.02	1978	0.05	0.08	2011	0.55	0.59
1883	-0.18	-0.19	1913	-0.38	-0.32	1946	-0.08	-0.02	1979	0.12	0.16	2012	0.57	*
1884	-0.26	-0.22	1914	-0.22	-0.30	1947	-0.04	-0.07	1980	0.23	0.15	2013	0.60	*
1885	-0.24	-0.25	1915	-0.16	-0.31	1948	-0.10	-0.10	1981	0.28	0.20	2014	*	*
1886	-0.23	-0.25	1916	-0.35	-0.29	1949	-0.11	-0.10	1982	0.09	0.20			
1887	-0.31	-0.21	1917	-0.43	-0.31	1950	-0.19	-0.09	1983	0.27	0.17			
1888	-0.19	-0.23	1918	-0.31	-0.33	1951	-0.06	-0.05	1984	0.12	0.14			
1889	-0.09	-0.23	1919	-0.28	-0.30	1952	0.02	-0.05	1985	0.08	0.18			
1890	-0.32	-0.23	1920	-0.26	-0.27	1953	0.09	-0.04	1986	0.14	0.19			
1891	-0.26	-0.26	1921	-0.20	-0.26	1954	-0.11	-0.06	1987	0.28	0.22			
1892	-0.30	-0.31	1922	-0.28	-0.25	1955	-0.12	-0.06	1988	0.35	0.28			
1893	-0.35	-0.29	1923	-0.25	-0.23	1956	-0.18	-0.07	1989	0.24	0.33			
1894	-0.32	-0.27	1924	-0.23	-0.21	1957	0.04	-0.04	1990	0.39	0.31			
1895	-0.24	-0.25	1925	-0.21	-0.19	1958	0.05	-0.02	1991	0.38	0.28			
1896	-0.17	-0.24	1926	-0.08	-0.17	1959	0.03	0.02	1992	0.19	0.29			
1897	-0.17	-0.21	1927	-0.17	-0.18	1960	-0.04	0.02	1993	0.21	0.30			
1898	-0.30	-0.19	1928	-0.16	-0.16	1961	0.05	0.03	1994	0.28	0.29			
1899	-0.19	-0.20	1929	-0.30	-0.16	1962	0.04	-0.01	1995	0.43	0.34			
1900	-0.14	-0.23	1930	-0.11	-0.15	1963	0.07	-0.03	1996	0.32	0.42			
1901	-0.20	-0.24	1931	-0.06	-0.16	1964	-0.20	-0.05	1997	0.45	0.44			
1902	-0.30	-0.28	1932	-0.10	-0.12	1965	-0.10	-0.06	1998	0.61	0.44			
1903	-0.36	-0.31	1933	-0.24	-0.13	1966	-0.04	-0.08	1999	0.40	0.48			
1904	-0.43	-0.32	1934	-0.09	-0.14	1967	-0.01	-0.03	2000	0.40	0.51			
1905	-0.29	-0.35	1935	-0.14	-0.11	1968	-0.05	-0.00	2001	0.52	0.51			
1906	-0.25	-0.36	1936	-0.10	-0.05	1969	0.06	-0.01	2002	0.61	0.53			
1907	-0.41	-0.37	1937	0.04	-0.03	1970	0.04	0.00	2003	0.60	0.58			
1908	-0.42	-0.40	1938	0.06	0.02	1971	-0.07	0.04	2004	0.51	0.59			
			1939	0.01	0.05	1972	0.02	0.02	2005	0.65	0.59			
			1940	0.07	0.06	1973	0.16	0.01	2006	0.59	0.57			
			1941	0.08	0.06	1974	-0.07	-0.01	2007	0.62	0.59			

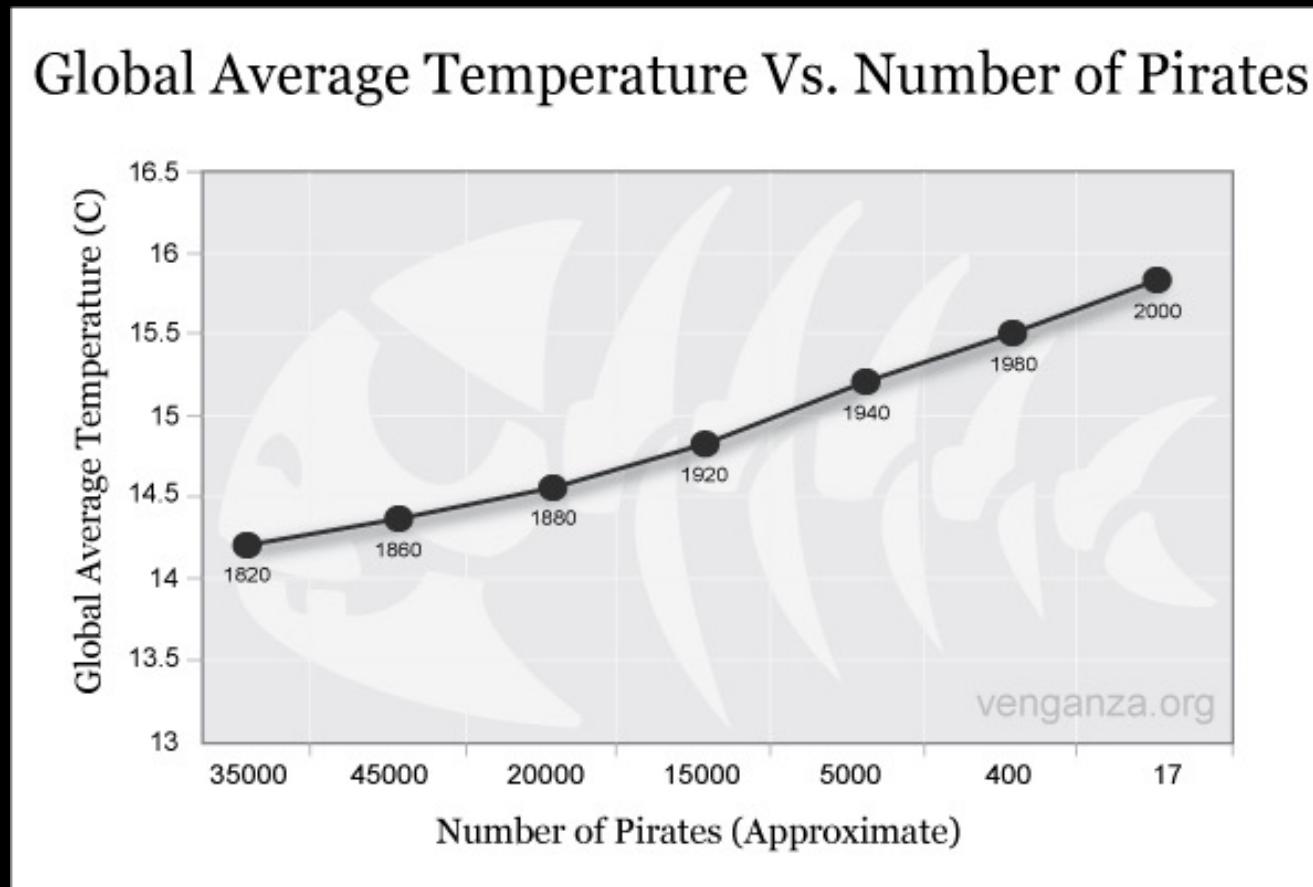
Global Means Temp as Graph



Viz Critical for Data Exploration



Statistics



Statistics: Nonsense Protection

I USED TO THINK
CORRELATION IMPLIED
CAUSATION.



THEN I TOOK A
STATISTICS CLASS.
NOW I DON'T.

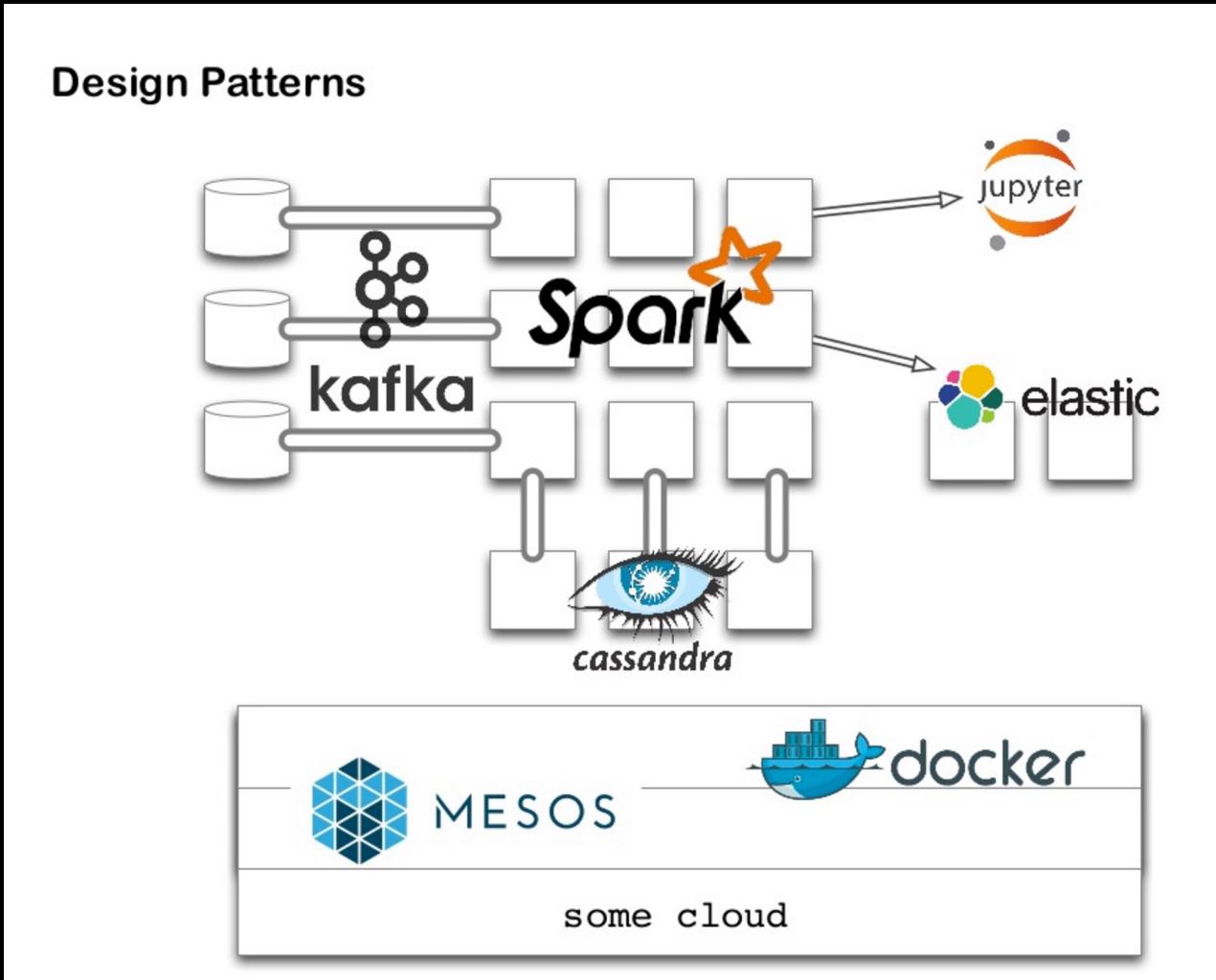


SOUNDS LIKE THE
CLASS HELPED.

| WELL, MAYBE.



Handling Big Data



Data Engineering

