ISE 599 Special Topics applied predictive Analytics

Mayank Kejriwal

Research Assistant Professor/Research Lead Department of Industrial and Systems Engineering Information Sciences Institute USC Viterbi School of Engineering

kejriwal@isi.edu

5-minute Review

Analytics: some broad views

- From Oxford Dictionary:
 - the systematic computational analysis of data or statistics
 - information resulting from the systematic analysis of data or statistics
- From Techopedia:
 - Analytics is the scientific process of discovering and communicating the meaningful patterns which can be found in data. It is concerned with turning raw data into insight for making better decisions. Analytics relies on the application of statistics, computer programming, and operations research in order to quantify and gain insight to the meanings of data. It is especially useful in areas which record a lot of data or information.

• From 'Innovating with Analytics', MIT Sloan Management Review:

 "...there is a strong correlation between driving competitive advantage and innovation with analytics and a company's effectiveness at managing the information transformation cycle, that is: capturing data, analyzing information, aggregating and integrating data, using insights to guide future strategy and disseminating information and insights."



Graphing Categorical Data

Categorical Data

(Qualitative Data)



Graphing Quantitative Data

Quantitative Data



Confounds

 Confounding occurs when the experimental controls do not allow the experimenter to reasonably eliminate plausible alternative explanations for an observed relationship between independent and dependent variables.

Correlation vs. Causation

 The use of a controlled study is the most effective way of establishing causality between variables. In a controlled study, the <u>sample</u> or <u>population</u> is split in two, with both groups being comparable in almost every way. The two groups then receive different treatments, and the outcomes of each group are assessed.

When making comparisons...

- ... Make sure to choose a suitable frame of reference!
- Single biggest cause of misleading (sometimes intentional) statistics, especially in business and science

Lesson 1: Statistics (+visualizations) are the first line of attack in most applied analytics pipelines "There are three kinds of lies – lies, damned lies, and statistics"

- To be an intelligent consumer of statistics, your first reflex must be to question the statistics you encounter
 - The more statistical knowledge and experience you have, the more adept you will be at doing the questioning
 - Like everything in life, skepticism should also be exercised in moderation
 - [Do not be like] Metrodorus of Chios: "None of us knows anything, not even this, whether we know or we do not know; nor do we know what 'to not know' or 'to know' are, nor on the whole, whether anything is or is not"

Takeaways

- Think about statistics like scales in music: no one really 'enjoys' it, but without mastering it, you cannot be a good musician
- In the classroom, we state the problem in an obvious way but in the real world modeling the problem is half the battle!

probability

What about probability?





	X
Notation	$\mathcal{N}(\mu,\sigma^2)$
Parameters	$\mu \in \mathbb{R}$ = mean (location)
	$\sigma^2 > 0$ = variance (squared scale)
Support	$x\in\mathbb{R}$
PDF	$rac{1}{\sqrt{2\pi\sigma^2}}e^{-rac{(x-\mu)^2}{2\sigma^2}}$
CDF	$\frac{1}{2} \left[1 + \operatorname{erf} \left(\frac{x - \mu}{\sigma \sqrt{2}} \right) \right]$
Quantile	$\mu+\sqrt{2}\mathrm{erf}^{-1}(2F-1)$
Mean	μ
Median	μ.
Mode	μ
Variance	σ^2
Skewness	0
Ex.	0
kurtosis	
Entropy	$rac{1}{2}\log(2\pi e\sigma^2)$
MGF	$\exp(\mu t + \sigma^2 t^2/2)$
CF	$\exp(i\mu t-\sigma^2 t^2/2)$
Fisher information	${\cal I}(\mu,\sigma) = egin{pmatrix} 1/\sigma^2 & 0 \ 0 & 2/\sigma^2 \end{pmatrix} {\cal I}(\mu,\sigma^2) = egin{pmatrix} 1/\sigma^2 & 0 \ 0 & 1/(2\sigma^4) \end{pmatrix} \end{pmatrix}$
Kullback- Leibler divergence	$D_{ ext{KL}}(\mathcal{N}_0 \ \mathcal{N}_1) = rac{1}{2} \{ (\sigma_0 / \sigma_1)^2 + rac{(\mu_1 - \mu_0)^2}{\sigma_1^2} - 1 + 2 \ln rac{\sigma_1}{\sigma_0} \}$

Let's play a real-estate parlor game

- There is a list of houses whose selling prices are available and given to you (training data)
 - Assume you have some statistical software/calculator available on hand and the data is big enough for statistical generalization
 - No other info (sq. ft., no. of bathrooms etc.) available other than selling prices
- Iteratively (say over 'ten' trials):
 - A house is randomly sampled from a new list (that is *similar* to the old list) whose selling prices are hidden from us
 - Everyone has to estimate the selling price of the house (if you want to do some statistics, you have time)
 - To simplify matters, pick one selling price as your estimate for all trials
 - After the trials are over, a final score is calculated and a winner is declared
- How do you **win** the game?

How is the score calculated?

- What's your estimate if the score is calculated as:
 - All or nothing (1 if you get the selling price exactly right, otherwise 0; doesn't matter if you're off by 1 USD or 1 million)?
 - Sum of deviations (sum over [true selling price estimate])?
 - Sum of **squared** deviations?

$$\gamma_{1} = \mathbf{E} \left[\left(\frac{X - \mu}{\sigma} \right)^{3} \right]$$
Intuition?
$$\begin{array}{l} \mathsf{Notation} \quad \mathcal{N}(\mu, \sigma^{2}) \\ \mu \in \mathbb{R} = \mathsf{mean}\left(\mathsf{location}\right) \\ \sigma^{2} > 0 = \mathsf{variance}\left(\mathsf{squared scale}\right) \\ \mathsf{Support} \quad x \in \mathbb{R} \\ \hline \mathsf{PDF} \quad \left[\frac{1}{\sqrt{2\pi\sigma^{2}}} e^{-\frac{(x - \mu)^{2}}{2\sigma^{2}}} \right] \\ \mathsf{CDF} \quad \left[\frac{1}{2} \left[1 + \operatorname{erf}\left(\frac{x - \mu}{\sigma\sqrt{2}} \right) \right] \\ \mathsf{Quantile} \quad \mu + \sigma\sqrt{2} \operatorname{erf}^{-1}(2F - 1) \\ \mathsf{Mean} \quad \mu \\ \mathsf{Modein} \quad \mu \\ \mathsf{Mode} \quad \varphi^{2} \\ \mathsf{Skewness} \quad 0 \\ \mathsf{Ex} \quad 0 \\ \mathsf{kurtosis} \\ \hline \mathsf{Entropy} \quad \left[\frac{1}{2} \log(2\pi e\sigma^{2}) \\ \mathsf{MGF} \quad \exp(\mu t + \sigma^{2} t^{2}/2) \\ \mathsf{CF} \quad \exp(\mu t - \sigma^{2} t^{2}/2) \\ \mathsf{Fisher} \\ \mathsf{Information} \quad \mathcal{I}(\mu, \sigma) = \left(\frac{1/\sigma^{2} - 0}{0} \right) \mathcal{I}(\mu, \sigma^{2}) = \left(\frac{1/\sigma^{2} - 0}{\sigma_{1}^{2}} - 1 + 2\ln \frac{\sigma_{1}}{\sigma_{0}} \right) \\ \mathsf{Kuilback} \\ \mathsf{eliver} \\ \mathsf{eli$$

х

$$\gamma_1 = \mathrm{E} iggl[iggl(rac{X-\mu}{\sigma} iggr)^3 iggr]$$

$$\operatorname{Kurt}[X] = \operatorname{E}\left[\left(rac{X-\mu}{\sigma}\right)^4\right]$$
 Intuition?

- Excess Kurtosis is K[x]-3
- Mesokurtic, Leptokurtic and platykurtic

	X
Notation	$\mathcal{N}(\mu,\sigma^2)$
Parameters	$\mu \in \mathbb{R}$ = mean (location)
	$\sigma^2 > 0$ = variance (squared scale)
Support	$x\in\mathbb{R}$
PDF	$rac{1}{\sqrt{2\pi\sigma^2}}e^{-rac{(x-\mu)^2}{2\sigma^2}}$
CDF	$rac{1}{2}\left[1+ ext{erf}igg(rac{x-\mu}{\sigma\sqrt{2}}igg) ight]$
Quantile	$\mu+\sigma\sqrt{2}\mathrm{erf}^{-1}(2F-1)$
Mean	μ
Median	μ
Mode	μ
Variance	σ^2
Skewness	0
Ex.	0
kurtosis	
Entropy	$\frac{1}{2}\log(2\pi e\sigma^2)$
MGF	$\exp(\mu t + \sigma^2 t^2/2)$
CF	$\exp(i\mu t-\sigma^2 t^2/2)$
Fisher information	${\mathcal I}(\mu,\sigma) = egin{pmatrix} 1/\sigma^2 & 0 \ 0 & 2/\sigma^2 \end{pmatrix} \end{pmatrix} {\mathcal I}(\mu,\sigma^2) = egin{pmatrix} 1/\sigma^2 & 0 \ 0 & 1/(2\sigma^4) \end{pmatrix} \end{pmatrix}$
Kullback- Leibler divergence	$D_{ ext{KL}}(\mathcal{N}_0 \ \mathcal{N}_1) = rac{1}{2} \{ (\sigma_0 / \sigma_1)^2 + rac{(\mu_1 - \mu_0)^2}{\sigma_1^2} - 1 + 2 \ln rac{\sigma_1}{\sigma_0} \}$
_	

$$\gamma_1 = \mathrm{E} iggl[iggl(rac{X-\mu}{\sigma} iggr)^3 iggr]$$

$$\operatorname{Kurt}[X] = \operatorname{E} \left[\left(rac{X-\mu}{\sigma}
ight)^4
ight]$$

$$H(X) = -\int_{-\infty}^{\infty} p(x) \log p(x) \, dx.$$

	Notation
	Parameters
	Support
	PDF
	CDF
	Quantile
	Mean
	Median
	Mode
	Variance
	Skewness
	Ex. kurtosis
	Entropy
	MGF
	CF
$\left(\sigma^{4}\right) $	Fisher information
$\ln \frac{\sigma_1}{\sigma_0} \}$	Kullback- Leibler divergence
σ	CF Fisher information Kullback- Leibler divergence

$$\gamma_1 = \mathrm{E} iggl[iggl(rac{X-\mu}{\sigma} iggr)^3 iggr]$$

$$\operatorname{Kurt}[X] = \operatorname{E}\left[\left(rac{X-\mu}{\sigma}
ight)^4
ight]$$

$$H(X) = -\int_{-\infty}^\infty p(x)\log p(x)\,dx.$$

The moment-generating function of a random variable X is

$$M_X(t):=\mathrm{E}ig[e^{tX}ig],\quad t\in\mathbb{R},$$

	X
Notation	$\mathcal{N}(\mu,\sigma^2)$
Parameters	$\mu \in \mathbb{R}$ = mean (location)
	$\sigma^2 > 0$ = variance (squared scale)
Support	$x\in\mathbb{R}$
PDF	$rac{1}{\sqrt{2\pi\sigma^2}}e^{-rac{(x-\mu)^2}{2\sigma^2}}$
CDF	$rac{1}{2}\left[1+ ext{erf}igg(rac{x-\mu}{\sigma\sqrt{2}}igg) ight]$
Quantile	$\mu+\sigma\sqrt{2}\mathrm{erf}^{-1}(2F-1)$
Mean	μ
Median	μ
Mode	μ
Variance	σ^2
Skewness	0
Ex.	0
kurtosis	
Entropy	$\frac{1}{2}\log(2\pi e\sigma^2)$
MGF	$\exp(\mu t + \sigma^2 t^2/2)$
CF	$\exp(i\mu t-\sigma^2 t^2/2)$
Fisher information	${\mathcal I}(\mu,\sigma) = egin{pmatrix} 1/\sigma^2 & 0 \ 0 & 2/\sigma^2 \end{pmatrix} {\mathcal I}(\mu,\sigma^2) = egin{pmatrix} 1/\sigma^2 & 0 \ 0 & 1/(2\sigma^4) \end{pmatrix} .$
Kullback- Leibler divergence	$D_{ ext{KL}}(\mathcal{N}_0 \ \mathcal{N}_1) = rac{1}{2} \{ (\sigma_0 / \sigma_1)^2 + rac{(\mu_1 - \mu_0)^2}{\sigma_1^2} - 1 + 2 \ln rac{\sigma_1}{\sigma_0} \}$

Statistical Tests

Null and Alternative Hypotheses

- Statistical hypothesis: claim about a parameter of a population.
- Null hypothesis (H₀): specifies a default course of action, preserves the status quo.
- Alternative hypothesis (H_a): contradicts the assertion of the null hypothesis,
- $\bullet H_a$: It is the Research Hypothesis \rightarrow
- > What you want to test
- > The question you want to investigate
- > The statement for which you are collecting evidence

Example problem

- The manager of a health maintenance organization has set as a target that the mean waiting time of non-emergency patients **will not exceed** 30 minutes (Status Quo). In spot checks, the manager finds the waiting times of 36 patients; the patients are selected randomly on different days. Assume that the population standard deviation of waiting times is 10 minutes, the sample mean is 35 minutes (What is given? Status Quo or Research Hypothesis?)
- a. What is the relevant parameter to be tested? $\mu = mean \ waiting \ time \ of \ ALL \ non-emergency$ patients
- b. Formulate null and research hypotheses.

Cont'd

c. State the test statistic and the rejection region corresponding to α = .05.

TestStatistic :

$$Z_{ts} = \frac{\overline{X} - \mu_0}{\sigma / \sqrt{n}}$$

$$Z_{ts} = \frac{35 - 30}{10 / \sqrt{36}} = \frac{5 * \sqrt{36}}{10} = 3.0$$
Rejection Region :
Reject Ho if Z > 1.645

$$Z_{ts} = 3.0$$

$$Z_{ts} = 3.0$$

Note 1: Central Limit Theorem (CLT)

- I will not state the CLT here, but I will say it is (typically) about the sampling distributions of sample means
 - Expected value of the sample mean does not change!
 - Variance declines over time, but it takes 30+ samples in practice before you can truly start assuming normality
- The theorem does not directly imply use of t-test etc. when doing significance testing (t-test assumes the data generator is normal to begin with, along with additional assumptions)
 - What do I do if data is non-normal? Use non-parametric tests!

Note 2: Estimator vs. parameter

- Parameters are **not** random variables (they're fixed and are the statistical equivalent of 'constants'; hence, we use greek letters)
- Estimators are almost always (I've never seen an exception but make no claims that there aren't) random variables, and can be point estimates or intervals
 - Terminology is overloaded, but estimators are generally understood to be 'functions' that yield an estimate (e.g., given a set of (x,y) points, what is an estimator for 'b' in the eqn. y=bx+a)
 - Estimators could be biased, inconsistent (if sample size can vary), inefficient, have high MSE (or all of the above)

➢ Use **bull's eye analogy** for intuition

- Rare to find efficient and consistent estimators, but they do exist, especially in linear least-squares regression!
- **Question:** What does it 'mean' when I say my 95% confidence interval for a sample mean is [0.3,0.6]?

Note 3: What if we have more than one null hypothesis?

- We can 'exchange' information between test statistics to gain power!
 - Global hypothesis testing: Simultaneously test all null hypotheses
 - Multiple hypothesis testing: separately test each null hypothesis

Global hypothesis testing

We assume we have *n* null hypotheses, $H_{0,i}$ with corresponding p-values p_i obtained as if we were only testing hypothesis *i* The global null hypothesis then is that all null hypotheses are true: $H_0 = \bigcap_{i=1}^n H_{0,i}$

Fisher's combination test:

- (1) Compute T=-2∑log(p_i) based on the individual p-values, noting that under the null hypothesis, a p-value has a uniform distribution on [0,1]
- (2) Under the global null hypothesis, T follows a chi-square distribution with 2n degrees of freedom, so we can decide to reject based on the value of T

Bonferroni method (better known in the engineering/science community)

Reject H_0 if $\min_i p_i \leq \alpha/n$

"The way to get started is to quit talking and begin doing." - Walt Disney