

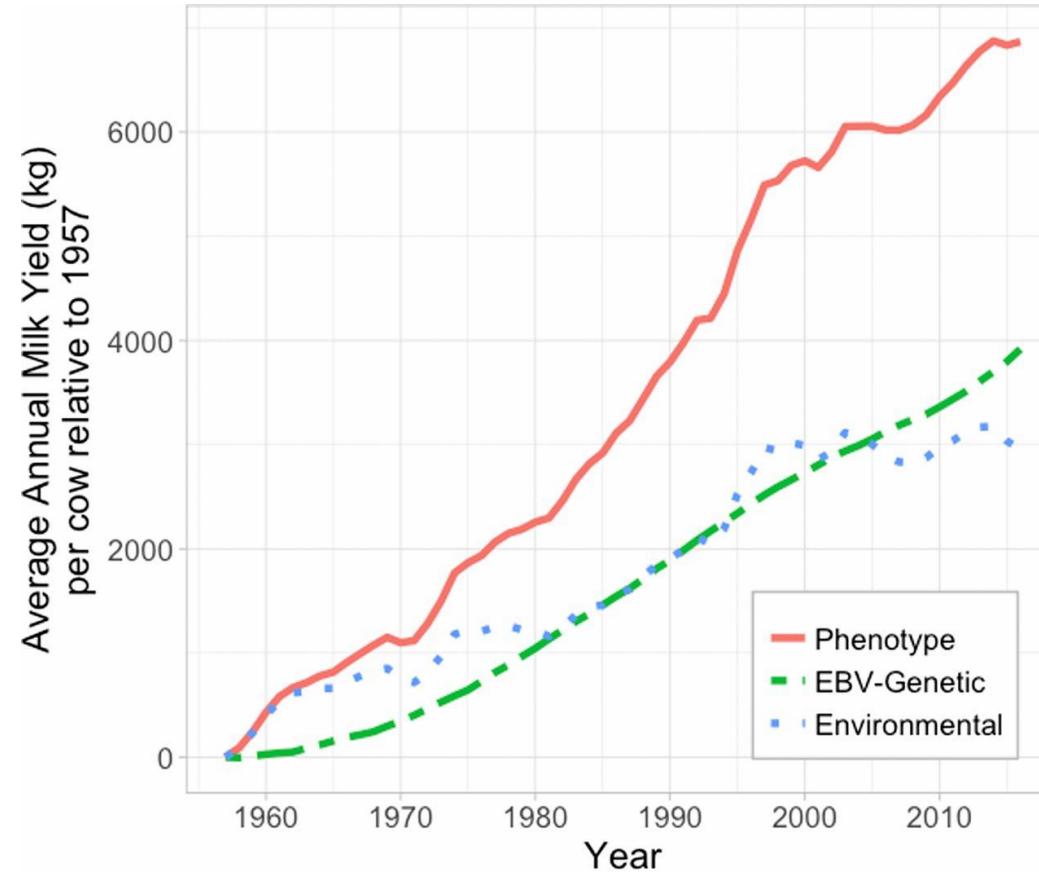


# Statistics and Data Science

Hao Cheng

Department of Animal Science  
University of California, Davis

# Why is Statistics important?



# Data Science

**#30 Bojangles'** 

**749 US Locations**



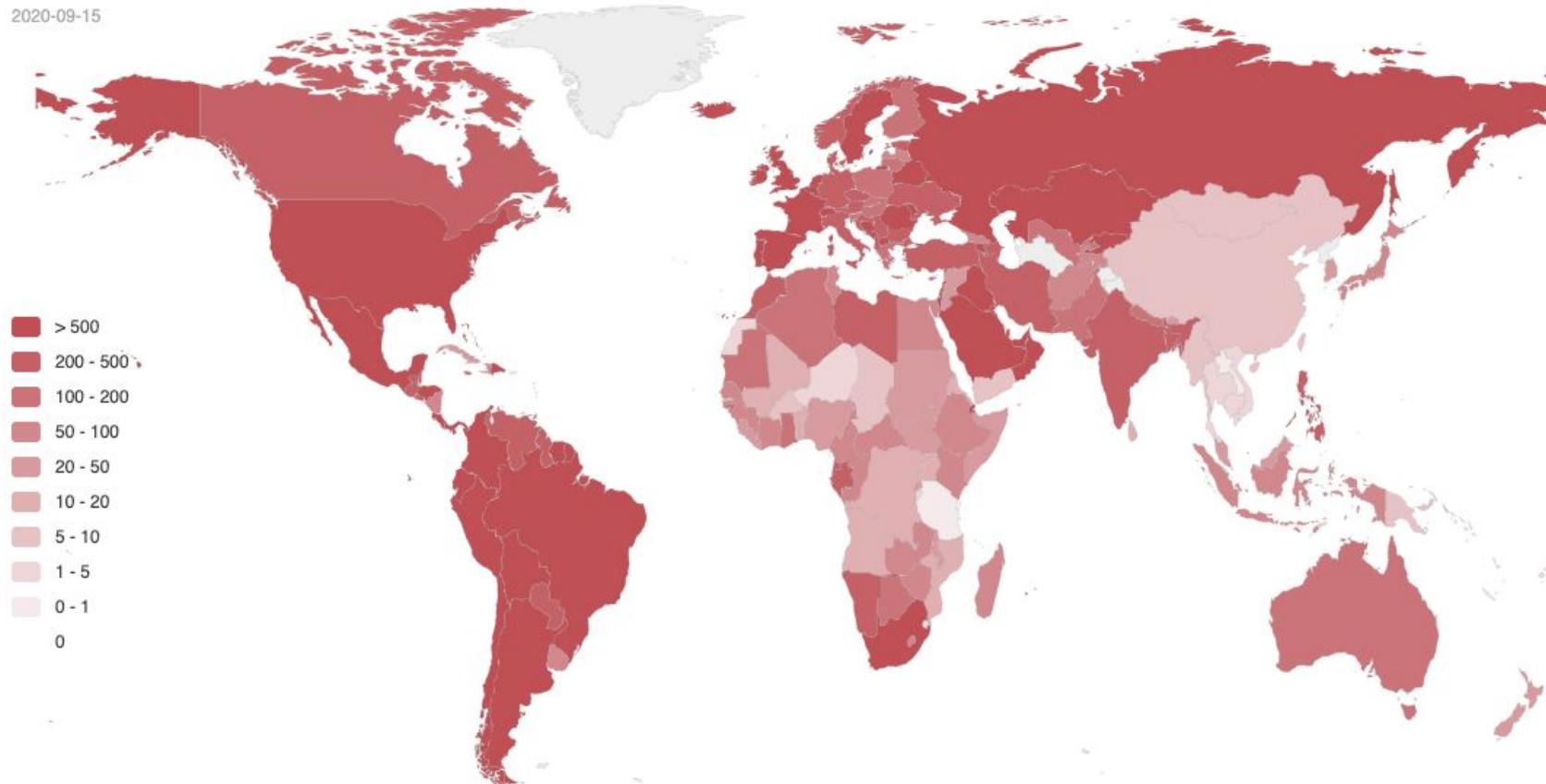
0:00 / 2:28



M. Anselme

## Number of Cases/100k Population

2020-09-15



# Pilot study: Application of artificial intelligence for detecting left atrial enlargement on canine thoracic radiographs

Shen Li, Zigui Wang, Lance C. Visser, Erik R. Wisner, Hao Cheng✉

First published: 11 August 2020 | <https://doi.org/10.1111/vru.12901>



# Quiz 1: What is p-value?

**The p-value is the probability of obtaining a value of the test statistic as or more extreme than the observed value of the test statistic when the null hypothesis is true.**

# Creative Writing Study

- Experimental units: writers
- Treatment: type of questionnaire given at the beginning of the study (questions emphasize either intrinsic or extrinsic rewards)
- Random assignment: 24 intrinsic, 23 extrinsic
- Response: average of 12 evaluations of creativity on a 40 point scale

## Observed data

<b>intrinsic:</b>	12.0	12.0	12.9	13.6	16.6	17.2
	17.5	18.2	19.1	19.3	19.8	20.3
	20.5	20.6	21.3	21.6	22.1	22.2
	22.6	23.1	24.0	24.3	26.7	29.7
 <b>extrinsic:</b>	5.0	5.4	6.1	10.9	11.8	12.0
	12.3	14.8	15.0	16.8	17.2	17.2
	17.4	17.5	18.5	18.7	18.7	19.2
	19.5	20.7	21.2	22.1	24.0	

null hypothesis:

Treatments have the same effect.

model-based

(relies on the specification of a model)

# model-based

- two independent random samples
- homogeneous population variances
- normality

# Randomization (design-based)

# Quiz 2: What is type I (II) error?

	True condition			
Total population	Condition positive	Condition negative	Prevalence $= \frac{\sum \text{Condition positive}}{\sum \text{Total population}}$	Accuracy (ACC) = $\frac{\sum \text{True positive} + \sum \text{True negative}}{\sum \text{Total population}}$
Predicted condition	True positive	False positive, Type I error	Positive predictive value (PPV), Precision = $\frac{\sum \text{True positive}}{\sum \text{Predicted condition positive}}$	False discovery rate (FDR) = $\frac{\sum \text{False positive}}{\sum \text{Predicted condition positive}}$
Predicted condition negative	False negative, Type II error	True negative	False omission rate (FOR) = $\frac{\sum \text{False negative}}{\sum \text{Predicted condition negative}}$	Negative predictive value (NPV) = $\frac{\sum \text{True negative}}{\sum \text{Predicted condition negative}}$
$\text{True positive rate (TPR), Recall, Sensitivity, probability of detection, Power}$ $= \frac{\sum \text{True positive}}{\sum \text{Condition positive}}$	$\text{False positive rate (FPR), Fall-out, probability of false alarm}$ $= \frac{\sum \text{False positive}}{\sum \text{Condition negative}}$	$\text{Positive likelihood ratio (LR+)} = \frac{\text{TPR}}{\text{FPR}}$	$\text{Diagnostic odds ratio (DOR)} = \frac{\text{LR+}}{\text{LR-}}$	$\text{F}_1 \text{ score} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$
	$\text{False negative rate (FNR), Miss rate}$ $= \frac{\sum \text{False negative}}{\sum \text{Condition positive}}$	$\text{Specificity (SPC), Selectivity, True negative rate (TNR)}$ $= \frac{\sum \text{True negative}}{\sum \text{Condition negative}}$	$\text{Negative likelihood ratio (LR-)} = \frac{\text{FNR}}{\text{TNR}}$	

# Common Statistical Tests

# Common statistical tests are linear models

Last updated: 28 June, 2019. Also check out the [Python version!](#)

See worked examples and more details at the accompanying notebook: <https://lindeloev.github.io/tests-as-linear>

Common name	Built-in function in R	Equivalent linear model in R	Exact?	The linear model in words	Icon
<b>Simple regression: <math>\text{Im}(y \sim 1 + x)</math></b>	<b>y is independent of x</b> P: One-sample t-test N: Wilcoxon signed-rank	<code>lm(y ~ 1)</code> <code>lm(signed_rank(y) ~ 1)</code>	✓ for N > 14	One number (intercept, i.e., the mean) predicts y. - (Same, but it predicts the <i>signed rank</i> of y.)	
	P: Paired-sample t-test N: Wilcoxon matched pairs	<code>lm(y1 - y2, paired=TRUE)</code> <code>wilcox.test(y1, y2, paired=TRUE)</code>	✓ for N > 14	One intercept predicts the pairwise $y_2 - y_1$ differences. - (Same, but it predicts the <i>signed rank</i> of $y_2 - y_1$ .)	
	<b>y ~ continuous x</b> P: Pearson correlation N: Spearman correlation	<code>cor.test(x, y, method='Pearson')</code> <code>cor.test(x, y, method='Spearman')</code>	✓ for N > 10	One intercept plus x multiplied by a number (slope) predicts y. - (Same, but with <i>ranked x</i> and y)	
	<b>y ~ discrete x</b> P: Two-sample t-test P: Welch's t-test N: Mann-Whitney U	<code>t.test(y1, y2, var.equal=TRUE)</code> <code>t.test(y1, y2, var.equal=FALSE)</code> <code>wilcox.test(y1, y2)</code>	✓ ✓ for N > 11	An intercept for <b>group 1</b> (plus a difference if <b>group 2</b> ) predicts y. - (Same, but with one variance <i>per group</i> instead of one common.) - (Same, but it predicts the <i>signed rank</i> of y.)	
<b>Multiple regression: <math>\text{Im}(y \sim 1 + x_1 + x_2 + \dots)</math></b>	P: One-way ANOVA N: Kruskal-Wallis	<code>aov(y ~ group)</code> <code>kruskal.test(y ~ group)</code>	✓ for N > 11	An intercept for <b>group 1</b> (plus a difference if group ≠ 1) predicts y. - (Same, but it predicts the <i>rank</i> of y.)	
	P: One-way ANCOVA	<code>aov(y ~ group + x)</code>	✓	- (Same, but plus a slope on x.) <i>Note: this is discrete AND continuous. ANCOVAs are ANOVAs with a continuous x.</i>	
	P: Two-way ANOVA	<code>aov(y ~ group * sex)</code>	✓	Interaction term: changing <b>sex</b> changes the y ~ <b>group</b> parameters. <i>Note: <math>G_{g \times s}</math> is an indicator (0 or 1) for each non-intercept levels of the group variable. Similarly for <math>S_{s \times k}</math> for sex. The first line (with G) is main effect of group, the second (with S) for sex and the third is the group × sex interaction. For two levels (e.g. male/female), line 2 would just be "S"; and line 3 would be S; multiplied with each G.</i>	[Coming]
	<b>Counts ~ discrete x</b> N: Chi-square test	<code>chisq.test(groupXsex_table)</code>	✓	Interaction term: (Same as Two-way ANOVA.) <i>Note: Run glm using the following arguments: <code>glm(model, family=poisson())</code>. As linear-model, the Chi-square test is <math>\log(y) = \log(N) + \log(\alpha) + \log(\beta) + \log(\alpha\beta)</math> where α and β are proportions. See more info in the accompanying notebook.</i>	Same as Two-way ANOVA
N: Goodness of fit	<code>chisq.test(y)</code>	<code>glm(y ~ 1 + G<sub>2</sub> + G<sub>3</sub> + ... + G<sub>N</sub>, family=...)^2</code>	✓	(Same as One-way ANOVA and see Chi-Square note.)	1W-ANOVA

List of common parametric (P) non-parametric (N) tests and equivalent linear models. The notation  $y \sim 1 + x$  is R shorthand for  $y = 1 + b + a \cdot x$  which most of us learned in school. Models in similar colors are highly similar, but really, notice how similar they *all* are across colors! For non-parametric models, the linear models are reasonable approximations for non-small sample sizes (see "Exact" column and click links to see simulations). Other less accurate approximations exist, e.g., Wilcoxon for the sign test and Goodness-of-fit for the binomial test. The signed rank function is `signed_rank = function(x) sign(x) * rank(abs(x))`. The variables G<sub>i</sub> and S<sub>j</sub> are "dummy coded" indicator variables (either 0 or 1) exploiting the fact that when  $\Delta x = 1$  between categories the difference equals the slope. Subscripts (e.g., G<sub>2</sub> or y<sub>i</sub>) indicate different columns in data. lm requires long-format data for all non-continuous models. All of this is exposed in greater detail and worked examples at <https://lindeloev.github.io/tests-as-linear>.

<sup>A</sup> See the note to the two-way ANOVA for explanation of the notation.

<sup>B</sup> Same model, but with one variance per group: `glm(value ~ 1 + G2, weights = varIdent(form = ~1|group), method="ML")`.



# Some Matrix Algebra

		movie			
		1	2	3	
customer	1	4	1	?	
	2	?	3	5	
	3	?	?	3	
	4	3	1	?	

Can we guess ratings  
for customer/movie  
combinations not  
in the dataset?

$y_{ij}$  = customer  $i$ 's rating  
of movie  $j$       Which movie is  
best?

$$y_{ij} = \mu + c_i + m_j + \epsilon_{ij}$$

		movie		
		1	2	3
customer	1	4	1	?
	2	?	3	5
	3	?	?	3
	4	3	1	?

Can we guess ratings  
for customer/movie  
combinations not  
in the dataset?

$y_{ij}$  = customer  $i$ 's rating  
of movie  $j$       Which movie is  
best?

$$y_{ij} = \mu_{ij} + \epsilon_{ij}$$