## Δμςδυσιρ γας



በበናና/Lላ<sup>6</sup> 1: ወዉ ህላ<sup>6</sup> ΓናበLCC<sup>6</sup> ርሲ ኦኒር / ፈካሪ ወቅር ወቅና, ወደ ኦ, ወደ ኦ, ውድር. የትርናበር ኦ/Lላ<sup>6</sup> ፈነት ህላ<sup>6</sup>: MODIS ኦሬ የኦጋ<sup>6</sup> ፈነትር ኦና/Lላ<sup>6</sup>, ላሮ 9, 2019 (NASA, 2019).

# SMAR省ICE

The communities in north Canada used to take advantage of the frozen fjords to travel from one place to another by snowmobiles.

Given a number of readings of ice thickness across time and space, could we recommend the people where/when it's safe to ride the snowmobile?



## Consequences



Credits: Beaverton Police Department

Here are marks for IA Algorithms questions last year:								
V N C	Women: Men: Dther:	[17, [18, [17,	14, 18, 18,	18, 11, 9,	12, 17, 9,	17, 17, 11,	] ] ]	
The mean marks are								
٨	Nomen:	13.22	(n=4	9)				
٦	Men:	12.28	(n=2	19)				
(	Other:	13.10	(n=1	0)				
Women do								

EXERCISE. How would you critique this analysis?

- This does not report confidence
- It's inappropriate to share this data or to report unaggregated data for scarcely represented categories
- It's drawing a general conclusion out of just one year of past data
  - On the other hand, if we just restrict ourselves to describe only what has happened already, and never say anything about the future, our ability to condition/shape the future would be restricted as well

Based on the model

Mark ~  $\mu_{\text{gender}} + N(0, \sigma^2)$ 

the 95% confidence intervals are

 $\hat{\mu}_F \in [11.8, 14.6]$  $\hat{\mu}_M \in [11.6, 12.9]$  $\hat{\mu}_O \in [10.0, 16.2]$ 

Women tend to do better than Men. There is too little data about Other to be confident in any comparison.

EXERCISE. How would you critique this revised analysis?

- Marks are not independent (each student answers to 2 questions)
- The Gaussian distribution does not seem appropriate
- If we want to report on differences, we should report a confidence interval for the differences

Based on a model using one-hot coding of gender,

Mark ~  $\mu_F + \delta_M 1_{\text{gender}=M} + \delta_O 1_{\text{gender}=O} + N(0, \sigma^2)$ 

the 95% confidence intervals are

 $\hat{\mu}_F \in [11.8, 14.6]$  $\hat{\delta}_M \in [-2.5, 0.6]$  $\hat{\delta}_{\Omega} \in [-3.6, 3.3]$ 

Neither  $\hat{\delta}_M$  nor  $\hat{\delta}_O$  is convincingly non-zero.

EXERCISE. How would you implement this analysis?

gender mark F 17 F 14 Μ 18 11 Μ Μ 17 : :

See Lecture 12---. # The readout function def t(marks): use sklearn.linear\_model to fit the proposed model to marks return a triple with the intercept ( $\mu_F$ ) and the coef ( $\delta_M, \delta_O$ ) Let  $\hat{\mu}_F, \hat{\delta}_M, \hat{\delta}_O, \hat{\sigma}$  be the mle estimates from the marks column in the dataset or, we could use def rmarks(): pred =  $\hat{\mu}_F + \hat{\delta}_M 1_{\text{gender}=M} + \hat{\delta}_O 1_{\text{gender}=O}$ return no random pormal (here in the dataset)

return np.random.normal(loc=**pred**, scale= $\hat{\sigma}$ )

# Get lots of samples of the test statistic t\_ = [t(rmarks()) for \_ in range(10000)] np.quantile([ $\theta$ [0] for  $\theta$  in t\_], [.025, .975]) # confint for  $\mu_F$  How might we decide whether this simpler model is good enough?

I think everyone gets pretty much the same mark, regardless of gender. Mark ~  $\mu$  + Normal(0,  $\sigma^2$ )

# To answer this, it can be helpful to introduce a richer model.



I think gender affects marks. Mark ~  $\mu_{\text{gender}}$  + Normal(0,  $\sigma^2$ )

#### confidence intervals

### FREQUENTIST

(The answer might depend on how we resample.)

### BAYESIANIST

(The answer depends on our priors for the unknowns.) For just two genders: Consider the richer model with  $\mu_{\text{gender}}$ and find a 95% confidence interval for  $\hat{\mu}_M - \hat{\mu}_F$ .

 $\mathbb{P}(\hat{\mu}_M - \hat{\mu}_F \in [-2.5, 0.6]) = 95\%$ This contains zero, i.e. I'm NOT confident that  $\hat{\mu}_M - \hat{\mu}_F$  is non-zero. So the simpler model is OK.

For just two genders: Consider the richer model with  $\mu_{\text{gender}}$ and find a 95% confidence interval for  $\mu_M - \mu_F$ .

 $\mathbb{P}(\mu_M - \mu_F \in [-3.1, -0.2]) = 95\%$ This does not contain zero, i.e. I am confident that  $\mu_M - \mu_F$  is non-zero. So the simpler model isn't good enough. If we have prior weights for two models (the simple model, and the richer model with  $\mu_{gender}$ ), we can find posterior weights using Bayes's rule.

For prior weights 50%/50%, the posterior weights are 79%/21% in favour of the simpler model.

This is great if there's a single model parameter that we want to investigate

This is for when we want to evaluate the model of a whole

#### **Hypothesis Testing**

# §9.3 HYPOTHESIS TESTING



Can you taste the difference between milk-first versus tea-first?

HYPOTHESIS: you can't.



### Fisher's hypothesis testing $t \leftarrow f \cdot t = 0$ Let x be the dataset.

State a null hypothesis  $H_0$ , i.e. a probability model for the dataset

- 1. Choose a test statistic  $t : \text{dataset} \mapsto \mathbb{R}$
- 2. Define a random synthetic dataset  $X^*$ , what we might see if  $H_0$  were true.
- 3. Look at the histogram of  $t(X^*)$ , and let p be the probability of seeing a value as extreme or more so than the observed t(x).

A low p-value is a sign that  $H_0$  should be rejected.

## x = taster's assignment of labels Ho: taster can't sell the difference. hence assighment is a random permutation of {t,t,t,t,m,m,m,m} t(x) = # correctdef X\*(): return random perm of {t,t,t,m,m,m,m} hist. & t(X\*) what world be the dist. of the fest statish?, if Ho were fine? $P = P(\epsilon(x^*) \neq \epsilon(x)) = 1.4$ p<5%: we'll reject the.

Degrees of freedom

- Probability model
- Null hypothesis
- Test statistic

Warning: assumptions!

#### Example 9.6.2.

I have a dataset with readings from two groups,  $x = [x_1, ..., x_m]$  and  $y = [y_1, ..., y_n]$ . Test whether the two groups are significantly different, using the test statistic  $\overline{y} - \overline{x}$ .

> 1 # 1. Define the test statistic 2 def t(x,y): return np.mean(y) - np.mean(x) 3 # 2. To generate a synthetic dataset, assuming H<sub>0</sub>, ... 4 xy = np.concatenate([x,y]) 5 def rxy\_star(): 6 return (np.random.choice(xy, size=len(x)), 7 np.random.choice(xy, size=len(y))) 8 # 3. Sample the test statistic under H0; find p-value for observed data 9 t\_ = np.array([t(\*rxy\_star()) for \_ in range(10000)]) 10 p = ...

#### Example 9.3.1.

I have a dataset with readings from two groups,  $x = [x_1, ..., x_m]$  and  $y = [y_1, ..., y_n]$ . Test whether the two groups are significantly different, using the test statistic  $\overline{y} - \overline{x}$ .

Ho: Xi, Yi both ~ N(M, J2) Equivalently, assume  $X_i \sim N(\mu, \sigma^2)$ ,  $Y_i \sim N(\mu + \delta, \sigma^2)$ H.: 8=0

1	<i># 1. Define the test statistic</i>
2	def t(x,y): return np.mean(y) - np.mean(x)
3	# 2. To generate a synthetic dataset, assuming $H_{\varrho}$ ,
4	<pre>xy = np.concatenate([x,y])</pre>
5	$\hat{\mu} = np.mean(xy)$
6	$\hat{\sigma}$ = np.sqrt(np.mean((xy - $\hat{\mu}$ )**2))
7	<pre>def rxy_star():</pre>
8	return (np.random.normal(loc= $\hat{\mu}$ , scale= $\hat{\sigma}$ , size=len(x)),
9	np.random.normal(loc= $\hat{\mu}$ , scale= $\hat{\sigma}$ , size=len(y)))
10	# 3. Sample the test statistic under H0; find p-value for observed data
11	<pre>t_ = np.array([t(*rxy_star()) for _ in range(10000)])</pre>

12  $p = 2 * \min(np.mean(t_ >= t(x,y)), np.mean(t_ <= t(x,y)))$ 

Two main questions

- What counts as 'more extreme'?
- How do we compute *p*?

## What counts as 'more extreme'?

- Plot the histogram for  $t(X^*)$ , assuming  $H_0$  is true
- Also plot the histogram for some scenarios where H<sub>0</sub> is false
- Do the alternatives push t(X\*) bigger, or smaller, or either? This determines what 'more extreme' means either one-tailed or two-tailed.

observed t  $H_0: \delta = 0$  $\delta > 0$  $\delta < 0$ -2 -4 0 2 4 if the observed t lies at either extreme, it's evidence against Ho: δ=0. How do we compute *p* for a two-tailed test?

The p-value is

$$\mathbb{P}\left(\begin{array}{c}t(X^*) \text{ at least}\\ \text{as extreme as }t(x)\end{array}\middle| H_0 \text{ is true}\right)$$



 $p = 2 * \min(np.mean(t_ >= t(x,y)), np.mean(t_ <= t(x,y)))$ 

Exercise 9.3.2 (Equality of group means).

We are given three groups of observations from three different systems,

$$x = [7.2, 7.3, 7.8, 8.2, 8.8, 9.5]$$
  

$$y = [8.3, 8.5, 9.2]$$
  

$$z = [7.4, 8.5, 9.0],$$

and we wish to know whether they all come from the same distribution, or whether there are three different distributions. Start with a general probability model in which they could potentially come from three different distributions,

$$X_i \sim \text{Normal}(a, \sigma^2), \quad Y_i \sim \text{Normal}(b, \sigma^2), \quad Z_i \sim \text{Normal}(c, \sigma^2)$$

Let  $H_0$  be that the three distributions are identical i.e. that a = b = c, or equivalently

$$X_i \sim Y_i \sim Z_i \sim \text{Normal}(\mu, \sigma^2).$$

Consider the test statistic

$$t = (\hat{a} - \hat{\mu})^2 + (\hat{b} - \hat{\mu})^2 + (\hat{c} - \hat{\mu})^2$$

where hats denote maximum likelihood estimators. If  $H_0$  were true, we'd expect t(x) to be small.

Find the value of the test statistic for the data given. What is the probability of seeing a value this large or larger, if  $H_0$  is true?

```
# 1. Define test statistic

def t(x,y,z):

\mu' = np.mean(np.concatenate([x,y,z]))

a',b',c' = [np.mean(v) \text{ for } v \text{ in } [x,y,z]]

return (a' - \mu')**2 + (b' - \mu')**2 + (c' - \mu')**2
```

```
#2. To generate a synthetic dataset, assuming H_0:
data = np.concatenate([x,y,z])
\hat{\mu} = np.mean(data)
\hat{\sigma} = np.sqrt(np.mean((data-\hat{\mu})**2))
```

```
def rxyz_star():
return (np.random.normal(size=len(x), loc=\hat{\mu}, scale=\hat{\sigma}),
np.random.normal(size=len(y), loc=\hat{\mu}, scale=\hat{\sigma}),
np.random.normal(size=len(z), loc=\hat{\mu}, scale=\hat{\sigma}))
```

```
#3. Sample the test statistic under H_0, find the p-value
t_sample = np.array([t(*rxyz_star()) for _ in range(10000)])
p = np.mean(t_sample >= t(x,y,z)) # answer: p = 0.592
```

The beauty of hypothesis testing is that it lets us test whether  $H_0$  is a good enough model for the data, without our having to specify an alternative model. Instead, we specify a test.

Where do test statistics come from?

There are two common scenarios, exploratory and rhetorical.

#### EXPLORATORY.

You, the modeller, are trying to come up with a good model for the dataset. Suppose you've tried out several models, and  $H_0$  is the best you've come up with. Is it good enough?

- If you settle for H<sub>0</sub> and someone else comes up with a better model, you lose.
- So it's up to you to creatively think up ways to test if H<sub>0</sub> might be deficient.

#### RHETORICAL.

Sometimes, there's a model  $H_1$  that everyone accepts to be the natural alternative to  $H_0$ .

- Example: H<sub>0</sub> = "my drug makes no difference", H<sub>1</sub> = "it makes a difference".
- If so, craft the test statistic to look for evidence pointing in the direction of H<sub>1</sub>.

### <sup>϶</sup>ͻϧϥϞͺͺͺͺͺͺͺͺ ͼϼϧϥͺͺͺͺͺͺͺͺ ͻϽϼϒͺͺͼϽͼϧϽϤϲ



በበናና/Lላ<sup>6</sup> 1: ወዉ ህላ<sup>6</sup> ΓናበLCC<sup>6</sup> ርሲ ኦኒር / ፈካሪ ወቅር ወቅና, ወደ ኦ, ወደ ኦ, ውድር. የትርናበር ኦ/Lላ<sup>6</sup> ፈነት ህላ<sup>6</sup>: MODIS ኦሬ የኦጋ<sup>6</sup> ፈነትር ኦና/Lላ<sup>6</sup>, ላሮ 9, 2019 (NASA, 2019).

# SMAR省ICE

The communities in north Canada used to take advantage of the frozen fjords to travel from one place to another by snowmobiles.

Given a number of readings of ice thickness across time and space, could we recommend the people where/when it's safe to ride the snowmobile?



# SMAR省ICE

The communities in north Canada used to take advantage of the frozen fjords to travel from one place to another by snowmobiles.

Given a number of readings of ice thickness across time and space: Model  $\leftarrow$  ice history t  $\leftarrow$  ice thickness H<sub>0</sub>  $\leftarrow$  "it's safe to ride in winter"

### **Iceberg detection**



Iceberg are detected by analysing SAR images. The backscatter of icebergs is very strong – very bright pixels. However, in rough waters there might be several lookalikes.

Given a number of readings of icebergs and ships in a given region, would we be able to tell the difference?



Subsets of Landsat image (left) and ASAR APP H-polarisation image (right) in the area where icebergs M, O and P were identified. Note that N was not observed in the Landsat image, but analysis of a nearby feature was done for the SAR images. These icebergs were not identified in the RADARSAT image. There are many more icebergs that can be identified in the Landsat image as well as in the SAR image. However, many of the bright points in the SAR image are not icebergs but speckle noise.

### **Iceberg detection**



Iceberg are detected by analysing SAR images. The backscatter of icebergs is very strong – very bright pixels. However, in rough waters there might be several lookalikes.

Given a number of readings of icebergs and ships in a given region, would we be able to tell the difference?  $H_0 \leftarrow$  end user won't be able to

tell the difference