# Experimental Design for Collection-based Comparative Evaluation

Evangelos Kanoulas
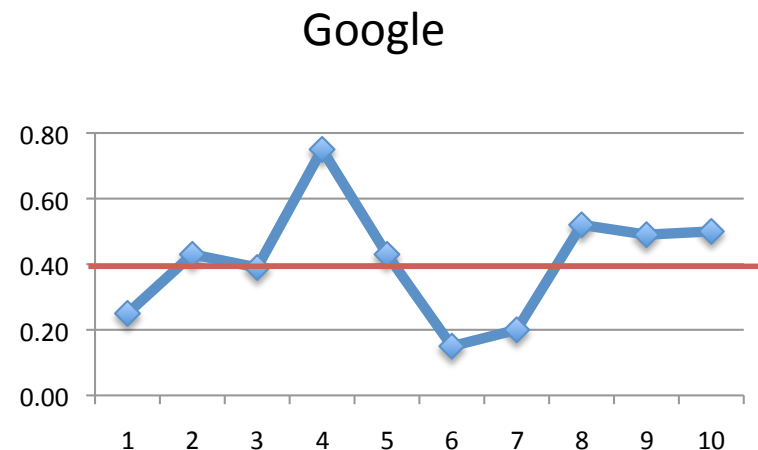
UNIVERSITY OF AMSTERDAM
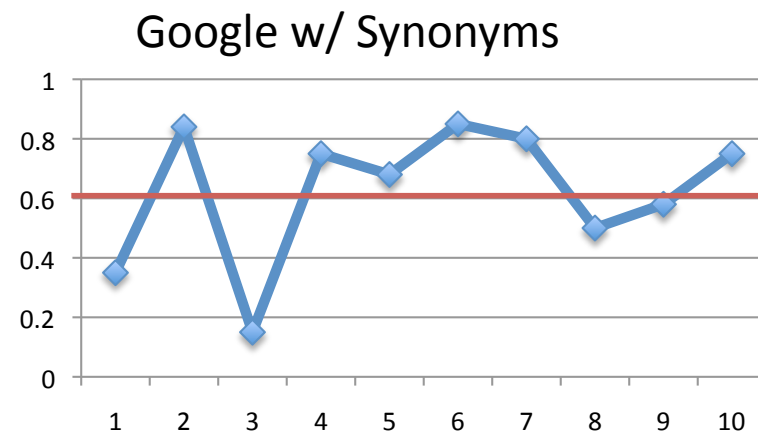
# Comparing Retrieval Systems

- Hypothesis: Synonyms will improve search engine effectiveness
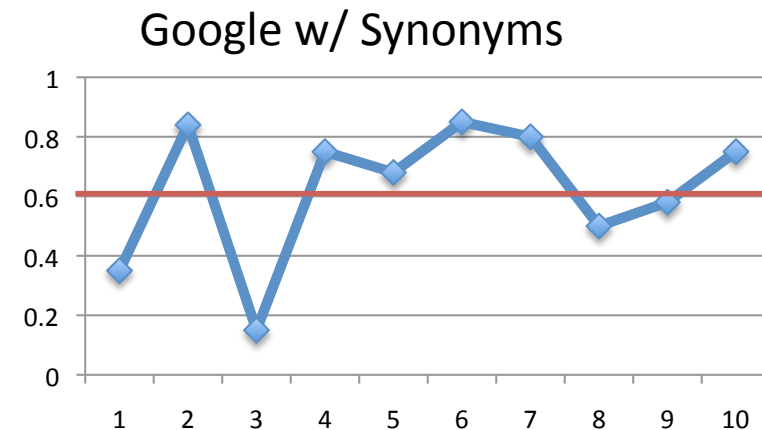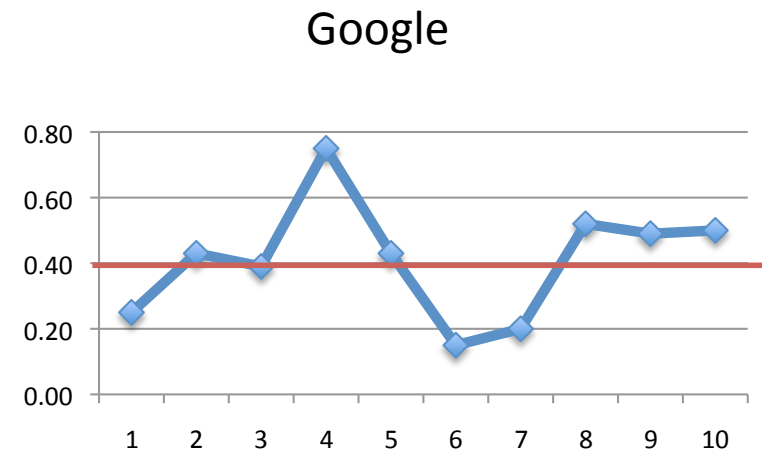
- Google:
  - Mean AP= 0.41

- Google
  w/ Synonyms
  - Mean AP= 0.63



Google



Google w/ Synonyms

# So what?

- Do these results support my hypothesis?

- Is it possible that my results are just random?

➡ statistical significance testing

### Google



### Google w/ Synonyms

# Statistical Significance Testing

- Two hypotheses, e.g.
  - $H_0$:  B-A = 0
  - $H_a$:  B-A ≠ 0 or B-A>0

- We want to prove the null hypothesis wrong

# Statistical Significance Testing

- Obtain system performance measurements over a sample of queries

- Compute a test statistic t from those measurements
  - with known distribution under $H_0$

- Compute the p-value, i.e.
  - the probability of observing the test statistic t …
  - … under a distribution obtained by assuming $H_0$ is true
    - If the p-value is low, conclude $H_0$ is false
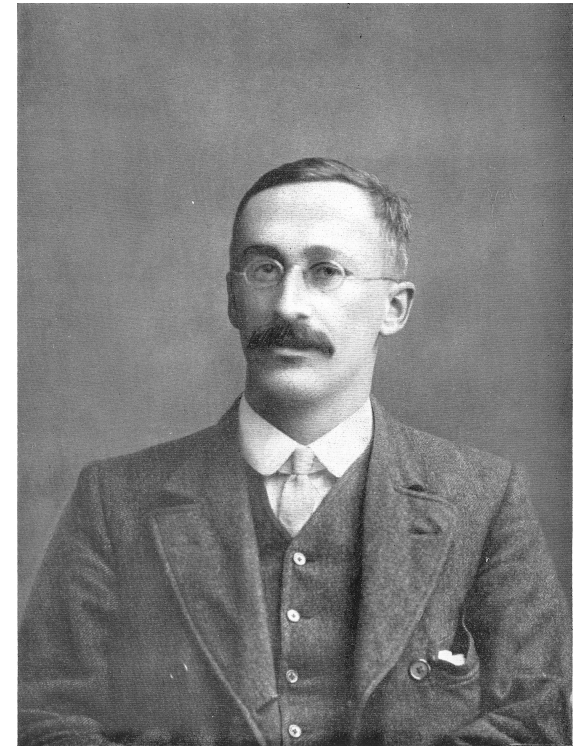
# Commonly used tests

- Non-parametric
  - Sign test/binomial test
  - Wilcoxon signed rank test

- Parametric
  - Student's t-test

- Distribution-free
  - Randomization test
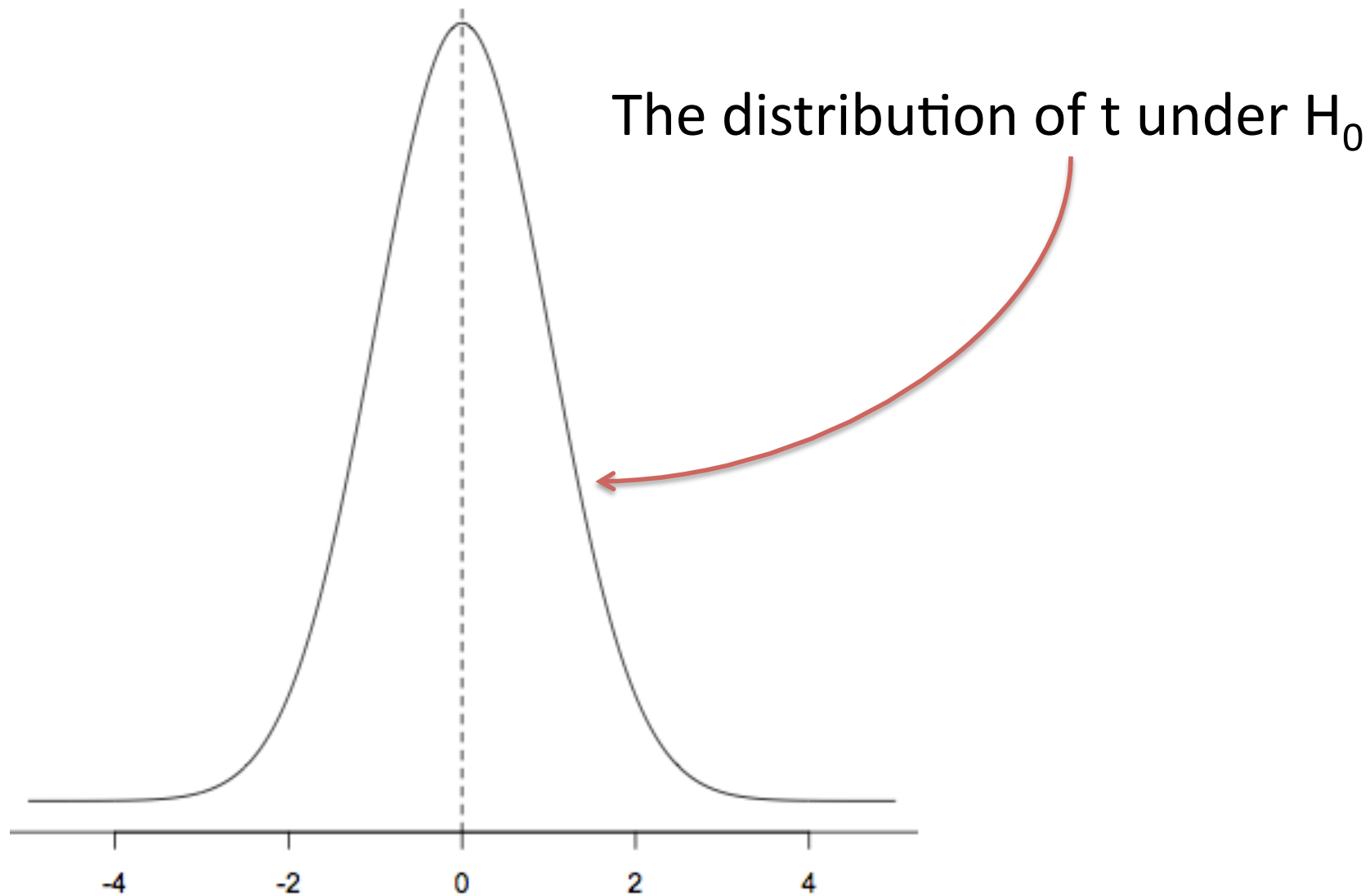  - Bootstrap test

# Student's t-test



- Statistic:
$$t = \frac{\overline{B - A}}{\frac{\sigma_{B-A}}{\sqrt{N}}}$$

   – (Assumption)
   - mean measures follow normal distribution

# Student's t-test



The distribution of t under $H_0$

# Student's t-test

| Query | A | B | B-A |
|-------|-----|-----|------|
| 1 | .25 | .35 | +.10 |
| 2 | .43 | .84 | +.41 |
| 3 | .39 | .15 | -.24 |
| 4 | .75 | .75 | 0 |
| 5 | .43 | .68 | +.25 |
| 6 | .15 | .85 | +.70 |
| 7 | .20 | .80 | +.60 |
| 8 | .52 | .50 | -.02 |
| 9 | .49 | .58 | +.09 |
| 10 | .50 | .75 | +.25 |

$$\hat{\mu} = \overline{B - A} = 0.214$$

$$\hat{\sigma}_{B-A} = 0.291$$

$$t = \frac{\hat{\mu}}{\hat{\sigma}_{B-A}} \sqrt{n} = 2.33$$

# Student's t-test

**Critical value (c.v.):**
Value of test statistic
when p = 0.05

Reject null hypothesis if:
- p value <= α

Or equivalently if
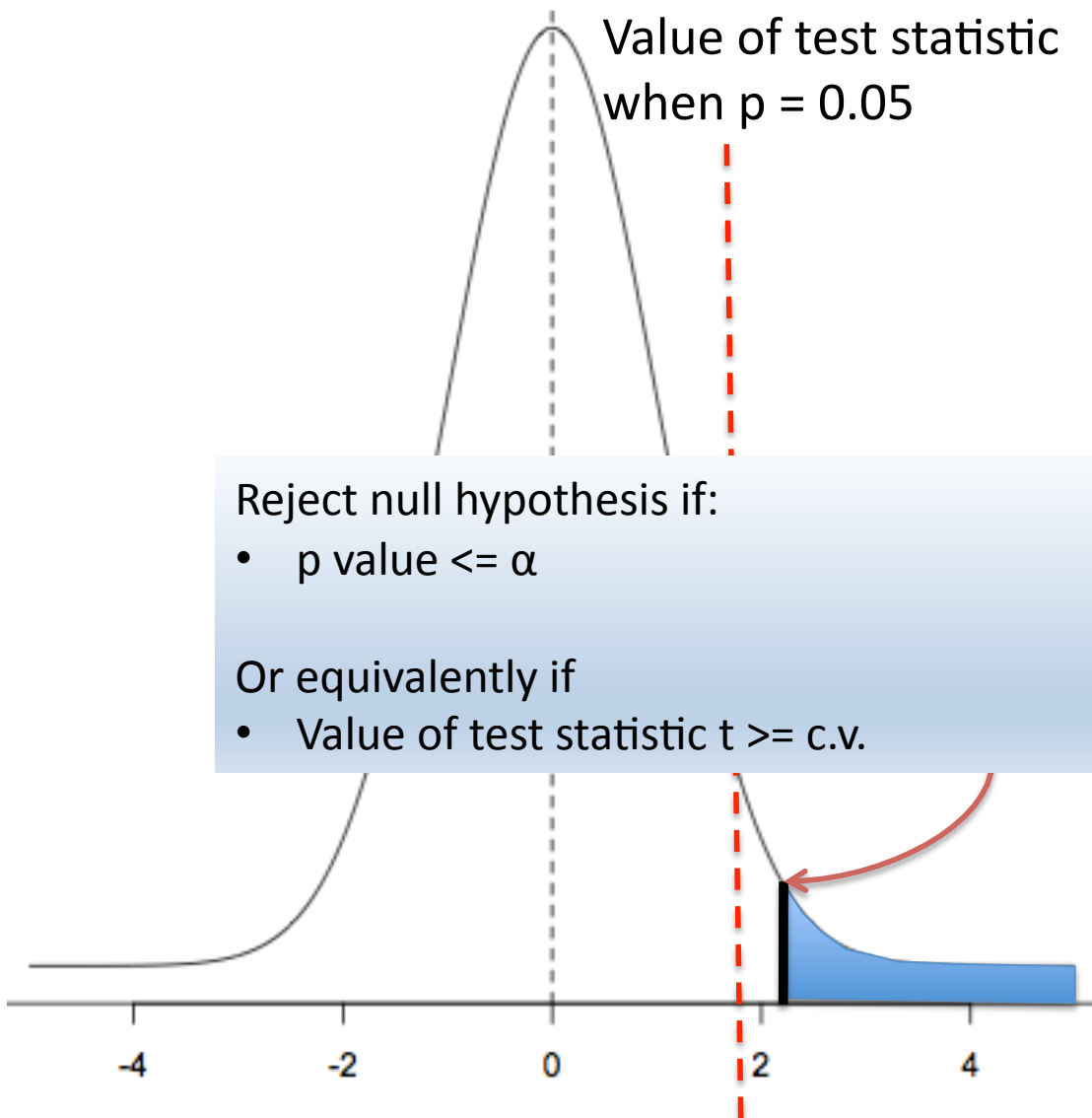- Value of test statistic t >= c.v.

$$\hat{\mu} = \overline{B - A} = 0.214$$
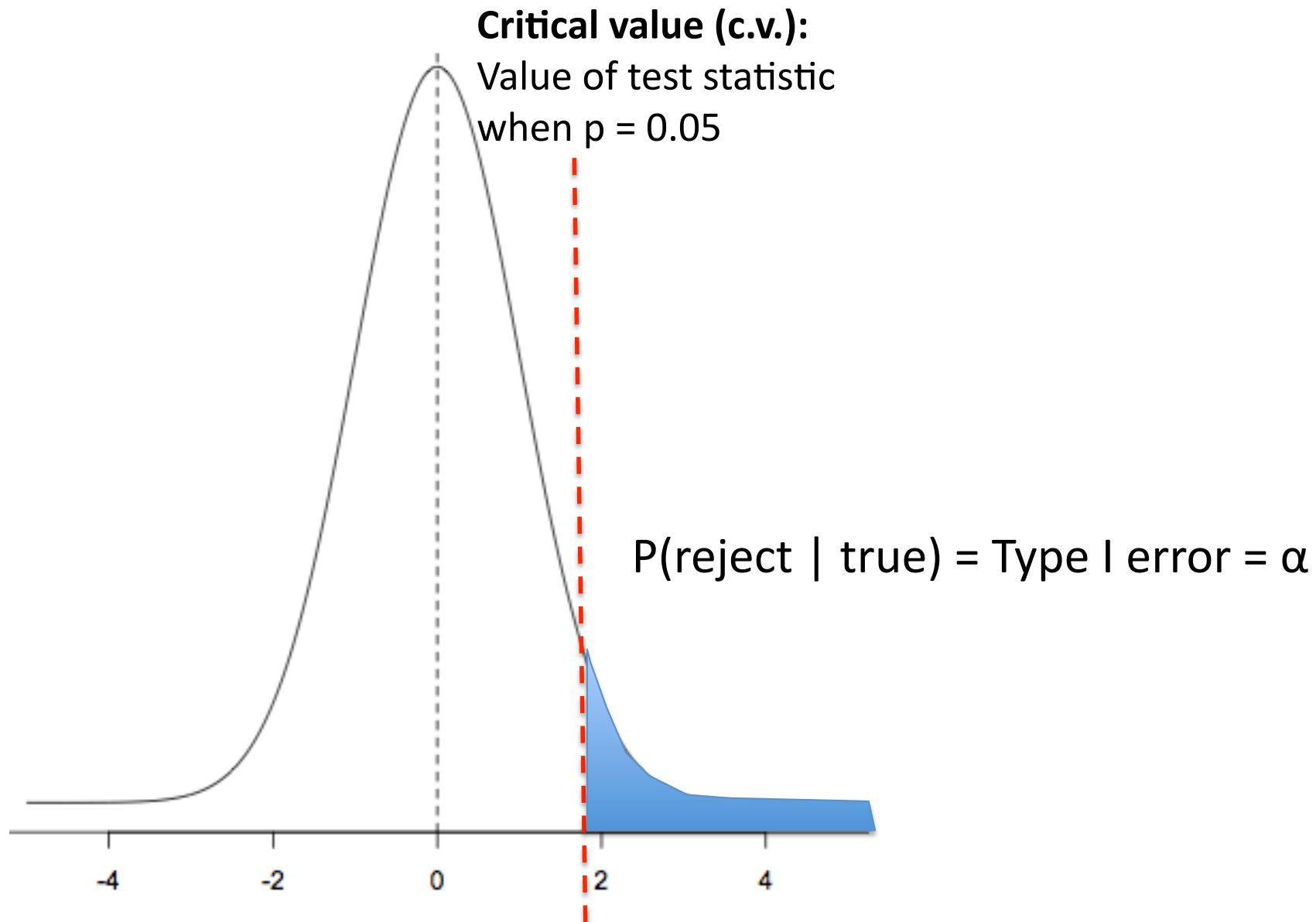
$$\hat{\sigma}_{B-A} = 0.291$$

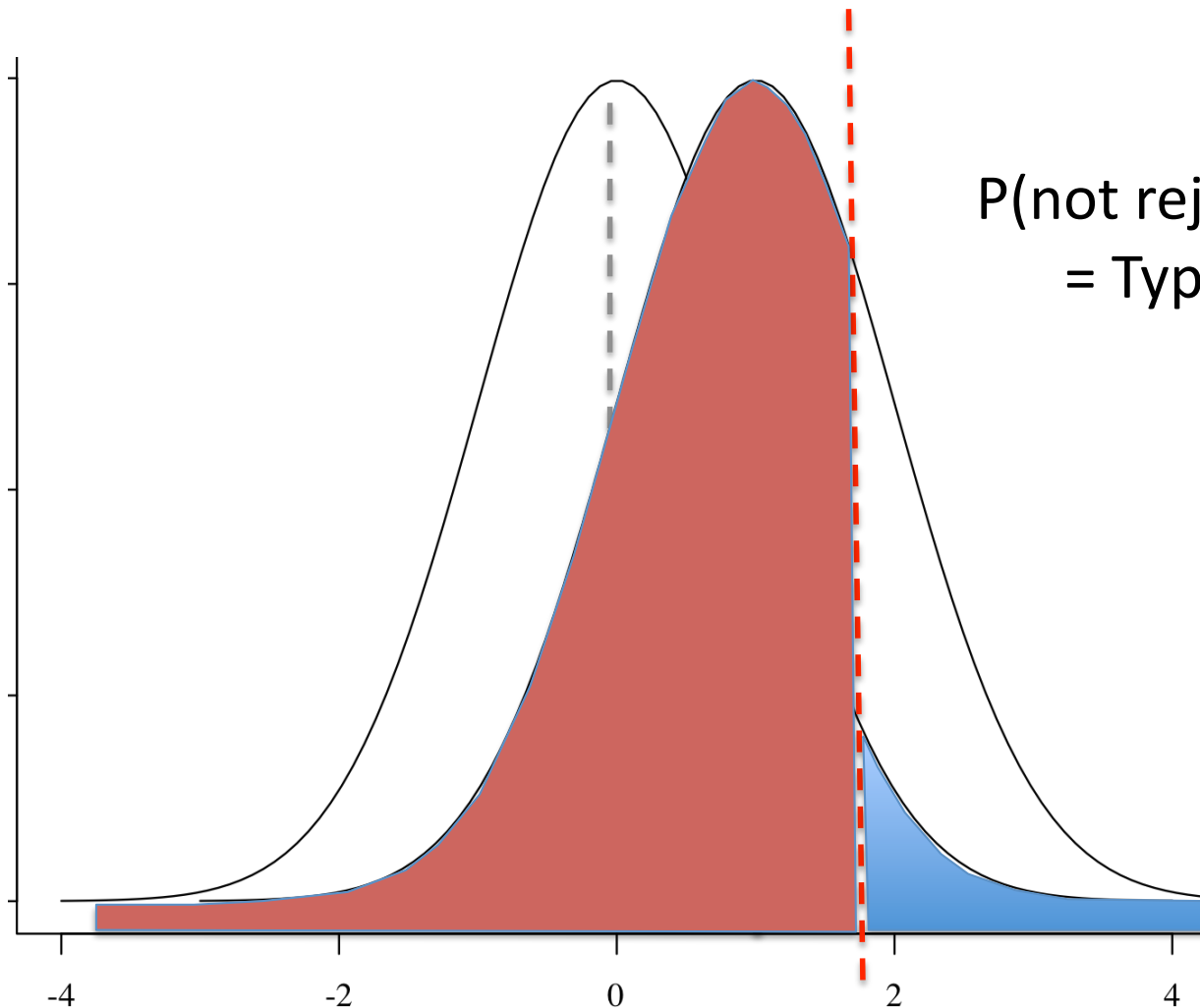$$t = \frac{\hat{\mu}}{\hat{\sigma}_{B-A}} \sqrt{n} = 2.33$$

$$p - value = 0.02$$

# Student's t-test



**Critical value (c.v.):** Value of test statistic when $p = 0.05$

$P(\text{reject} \mid \text{true}) = \text{Type I error} = \alpha$

# Student's t-test



**Critical value (c.v.):**
Value of test statistic
when p = 0.05

P(not reject | false)
    = Type II error = β

# Errors in Inference

- A significance test is basically a classifier

| H$_0$ | true | false |
|---|---|---|
| **not rejected** | accuracy: 1-α | Type II error: **β** |
| **rejected** | Type I error: **α** | power: 1-β |

- We can't actually know whether H$_0$ is true or not
  - If we could, we wouldn't need the test

- But we set up the test to control the expected Type I (significance) and Type II (power) error rates

# Expected Type I Error Rate

- Test parameter $\alpha$ is used to decide whether to reject $H_0$ or not—if $p < \alpha$, then reject $H_0$

- Choosing $\alpha$ is equivalent to stating an expected Type I error rate
  - e.g. if $p < 0.05$ is considered significant, we are saying that we expect that we will incorrectly reject $H_0$ 5% of the time

- Why?
  - Because when $H_0$ is true, every p-value is equally likely to be observed
  - 5% of the time we will observe a p-value less than 0.05... and therefore there is a 5% Type I error rate

# Expected Type II Error Rate

- What about Type II errors?
  - False negatives are bad: if we can't reject $H_0$ when it's false, we may miss out on interesting results

- What is the distribution of p-values when $H_0$ is false?
  - Problem: there is only one way $H_0$ can be true, but there are many ways it can be false

# Student's t-test



Effect Size

$$t = \frac{\overline{B - A}}{\frac{\sigma_{B-A}}{\sqrt{N}}}$$

# Effect Size

- A measure of the magnitude of the difference between two systems

  – Effect size is dimensionless; intuitively similar to % change in performance

  – Bigger population effect size => more likely to find a significant difference in a sample

# Power and Effect Size

- Before testing, we can say "I want to be able to detect an effect size of h with probability β"

  – "If there is at least a 5% difference, the test should say the difference is significant with 80% probability"

  $\Rightarrow$ h = 0.05, β = 0.8

# Sample Size

- Once we have chosen α, β, h, we can determine the sample size needed to make the error rates come out as desired
  - $n = f(α, β, h)$
  - Usually involves a linear search
  - There are software tools to do this

- Basically:
  - Sample size n increases with β if other parameters held constant
  - If you want more power, you need more queries

# Sample Size



$$t = \dfrac{\overline{B - A}}{\dfrac{\sigma_{B-A}}{\sqrt{N}}}$$

# Power Analysis

- Statistical significance testing:
  1. sample size
  2. effect size = diff of means / st. dev.
  3. significance level = P(Type I error) = probability of finding an effect that is not there
  4. power = 1 - P(Type II error) = probability of finding an effect that is there

- Given any three, we can determine the fourth
  - Easier under normality assumption

# So far... statistics 101...

- Two sides of the same coin:

  - Statistical significance => results generalize from a sample of queries to the population

  - Power analysis => number of queries necessary to stat. detect a given difference

# Why care about significance testing?

- Sources of variance specific to IR:
  - Properties of queries
  - Properties of document corpus
  - Properties of effectiveness measures
  - Assessor error and disagreement
  - Missing relevance judgments
  - Total number of relevant documents
  - ...

- Only variance due to queries included in standard statistical testing

  => Wrong conclusions!

# Collection-based Experiment
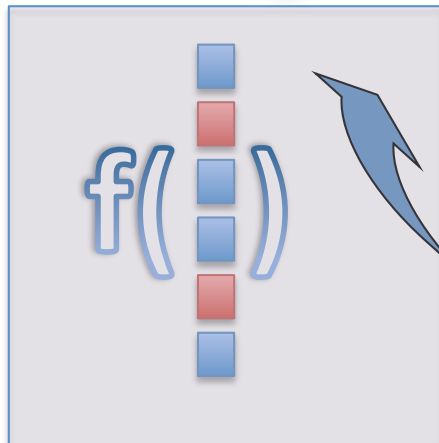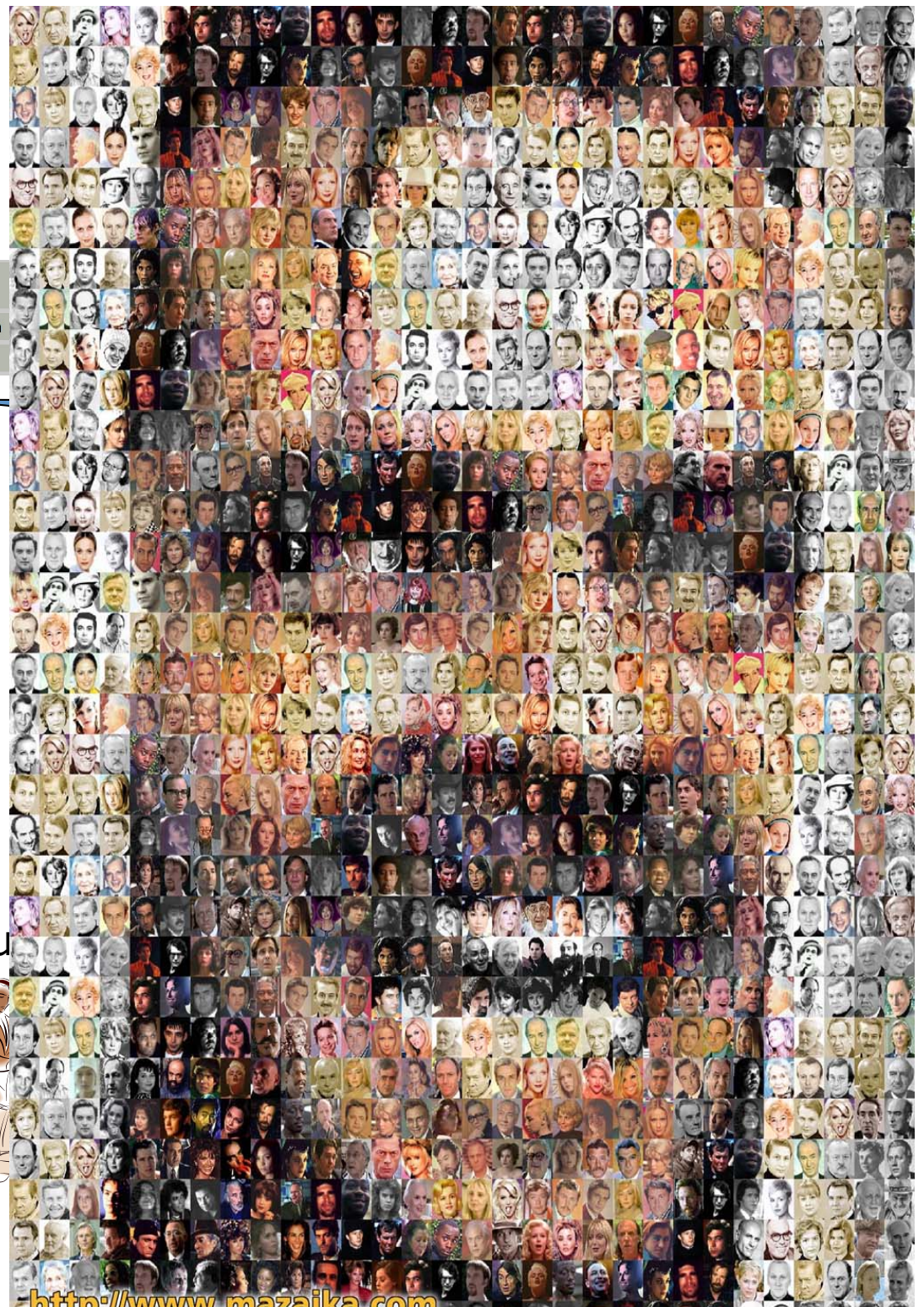
# Variance due to Queries



Results

Search Algorithms

# Variance due to Document Collection

uses of alternative dispute resolution

○ job search vancouver washington

○ poem of arrival of columbus

Results

Search Algorithms

# Variance due to missing judgments

$f(\quad)$

Judges

Ticino. Valais. Zürich Region. According to grade of difficulty: Conditions: .... The map also shows the entire hiking trail network (without additional information) ...

Ticino / Tessin : Climbing, Hiking & Mountaineering : SummitPost
www.summitpost.org › Alps - Western Part › Lepontine Alps ▾
Nov 9, 2007 – Ticino / Tessin : SummitPost.org : Climbing, hiking, mountaineering. ... real beautiful hiking trails and challenging peaks waiting to be climbed.

Hiking Switzerland, hiking trails for hikers, holiday resorts in ...
www.switzerland-hiking.ch/ ▾
The website for hikers and nature-lovers with information on hiking trails and ... Swiss panorama - Hiking trails and walking holidays .... Hiking Trails Ticino ...

# Variance due to Measure Parameters

uses of alternative dispute resolution

○ job search vancouver washington

○ poem of arrival of columbus

### Search Algorithms

Yandex

Bai 度

bing Beta

Google

f( )

Ju

http://www.mazaika.com

# Variance due to Document Collection



uses of alternative dispute resolution

job search vancouver washington

poem of arrival of columbus

Results

Search Algorithms

# Variance due to Document Collection

- The document collection is not absolute
  - may think of it as a sample
    - from some large/infinite universe of possible items

- Each query measurement is an estimate
  - of a population measure
    - one query, population of documents

- Quality of estimate varies between topics
  - therefore a mean is misleading
    - and so is a t-test

# Motivation Hypothesis

A test document collection
should be thought of as
a sample from some hypothetical universe
of possible documents

# Statistical significance

- Traditional significance testing:
  - consider the *queries* as a sample from some universe
    - what does this sample tell us about the population?

A test document collection should be thought of as a sample from some hypothetical universe

- Now we have two simultaneous sampling processes
  - need to revise the question
    - what does this (sample x sample) tell us about the (population x population)?

# Simulation of multiple collections



- How can you compute a measure over multiple collections?
  - Consider multiple collections
  - Simulate multiple collections

# Single Query Measurements

# Single System Measurements

# The Linear Model

- The t-test is based on a linear regression model

# The Linear Model

- The t-test is based on a linear regression model

the value of a measure calculated
on query *j* for system *i*

$\beta_i$ : system *i*

$$y_{ij} = \beta_i + b_j + \epsilon_{ij}$$

b$_j$ : query j

$\epsilon_{ij}$ : residual error

$$b_j \sim \mathcal{N}(0, \sigma_1^2)$$

$$\epsilon_{ij} \sim \mathcal{N}(0, \sigma^2)$$

# The Linear Model

- In the statistical programming environment **R**

```
lme(effectiveness ~ system, data=data, random=~1|
    query)
```

... equivalent to ...

```
t.test(effectiveness ~ system, data=data,
paired=TRUE)
```

# Mixed Effects Models

- Two sources of variance
  - Query effect
  - Collection effect (within query variance)

the value of a measure calculated
on query *j* for system *i*
and collection k

$b_j$ : query j

$\varepsilon_{ijk}$ : residual error

$$y_{ijk} = \beta_i + b_j + c_{ij} + \epsilon_{ijk}$$

$\beta_i$ : system *i*

$c_{ij}$: within query j effect

$$b_j \sim \mathcal{N}(0, \sigma_1^2)$$

$$c_{ij} \sim \mathcal{N}(0, \sigma_2^2)$$

$$\epsilon_{ijk} \sim \mathcal{N}(0, \sigma^2)$$

# Mixed Effects Models

- In the statistical programming environment **R**

```
lme1 <- lme(effectiveness~system, data=df, random=~1|query/system)

summary(lme1)

Random effects:
  Formula: ~1 | query
          (Intercept)
StdDev:   1.539644

  Formula: ~1 | system %in% query
          (Intercept)   Residual
StdDev:       0.6191864 0.6386645

Fixed effects: y ~ system
                  Value    Std.Error      DF      t-value        p-value
(Intercept)      -1.3445  0.2438470      846      -5.514077       0.0000
system2           0.0999  0.1343512      46        0.744112       0.4606
```

Mixed-effects Homoscedastic Model

# Mixed Effects Models

- Two sources of variance
  - Query effect
  - Collection effect (within query variance)

$$y_{ijk} = \beta_i + b_j + c_{ij} + \epsilon_{ijk}$$

$$b_j \sim \mathcal{N}(0, \sigma_1^2) \qquad \boxed{c_{ij} \sim \mathcal{N}(0, \sigma_{ij}^2)} \qquad \epsilon_{ijk} \sim \mathcal{N}(0, \sigma^2)$$

Heteroscedastic Model

Mixed-effects Heteroscedastic Model

# Mixed Effects Models

- In the statistical programming environment **R**

```
lme2 <- lme(effectiveness~system, data=df, random=~1|query/system,
    weights=varIdent(form=~1|query*system))
```

```
Random effects:
  Formula: ~1 | query              Formula: ~1 | system %in% query
        (Intercept)                       (Intercept)    Residual
StdDev:  1.447164                  StdDev:   0.4537618   0.186183
```

```
Variance function:
  Structure: Different standard deviations per stratum
  Formula: ~1 | query* system
  Parameter estimates:
        1*1       1*2          2*1          2*2             ...
        1.0000000 1.6108387    1.3969085    1.5405710       ...
```

```
Fixed effects: y ~ system
                Value     Std.Error      DF    t-value      p-value
  (Intercept)        -1.4385   0.22266286   846  -6.460817    0.0000
  system2        0.1834     0.09844907    46    1.863342       0.0688
```
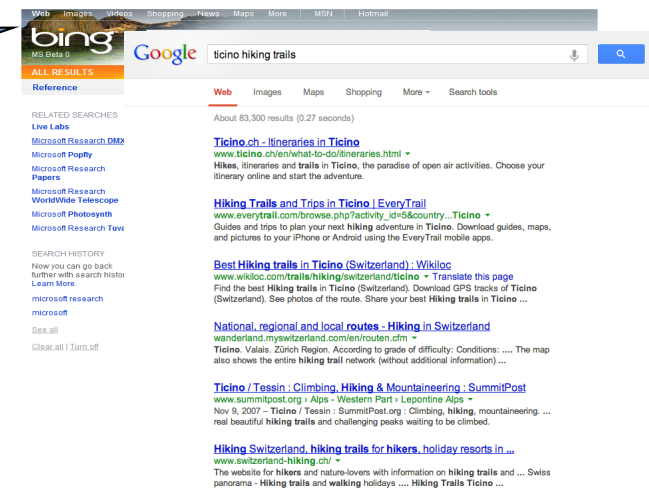
# Variance due to Measure Parameters

uses of alternative dispute resolution

○ job search vancouver washington

○ poem of arrival of columbus

### Search Algorithms

Yandex
Baidu百度
bing
Google

$f(\ |\ )$

Ju...

http://www.mazaika.com

# RBP – User Model

# RBP – The Measure

1

2

3

4

5

6

7

8

9

10

...

$$RBP = \sum_{i=1}^{n} rel_i (1-\vartheta)^{i-1} \vartheta$$

Relevance discounted by geometric distribution

# Choosing Parameter Values

- Different approaches:
  - Predefine parameters
  - Use click log; fit a model to gaps between clicks (Zhang et al., IRJ, 2010)
  - Minimize variance in evaluation (Kanoulas & Aslam, CIKM '09)

- All user models have parameters
  - Metrics evaluated at fixed parameter values
  - Evaluation w.r.t. an average user

# Choosing Parameter Values

- Users behave very differently when they search
  - Distribution of parameters (users) need to be considered

- A different approach
  - Mine Web Query logs
  - Learn a distribution of the parameters
  - Use this distribution to evaluate the quality of systems

# Patience Distribution for RBP

- Goal: produce a posterior distribution for θ
- Start with a uniform distribution for θ
- Update it based on logged data

# Posterior Distribution of Patience θ for RBP

$$P(\theta \mid E) = P(\theta \mid c) = \sum_{r=0}^{\infty} P(\theta \mid r,c)P(r \mid c)$$

$$P(\theta \mid r,c) \propto P(c \mid \theta,r)P(\theta \mid r)$$

The probability that user skips r document

$$P(c \mid \theta,r) = NB(r,\theta)$$

$$P(\theta \mid r) = Beta(\alpha,\beta)$$

- Start with uniform prior (α=β=1)

Probability distribution of the number of successes in a sequence of Bernoulli trials before **r** failures occurs.

# Posterior Distribution of Patience θ for RBP

$$P(\theta \mid r, E) \propto P(E \mid \theta, r) P(\theta \mid r)$$

$$P(E \mid \theta, r) = NB(r, \theta) \qquad P(\theta \mid r) = Beta(\alpha, \beta)$$

- If there are m queries, with r number of failures, and $c_i$ number of successes, i=1..m

$$P(\theta \mid r, E) = Beta(\alpha + \sum_{i=1}^{m} c_i, \beta + mr)$$

# Posterior Distribution for Impatience: RBP

- Distribution of users using the AOL log

# Distribution of RBP

- RBP values for different users for a single system and a single query

# Mixed Effects Models

$$y_{ijk} = \alpha_i + (\beta_j + \phi_j p_k) + (\kappa_{ij} + \gamma_{ij} p_k) + \epsilon_{ijk}$$

$y_{ijk}$: value of a metric on topic j for system i with parameter $p_k$

$\alpha_i$ : effect of system i

$\beta_j$ : effect of topic j

$\varphi_j p_k$ : interaction of topic with RBP parameter

$\kappa_{ij}$ : system/topic interaction effect

$\gamma_{ij} p_k$ : interaction of system/topic with RBP parameter

$\varepsilon_{ijk}$ : system/topic/parameter interaction effect

**TB06 NP**

# Variance due to Query Intents



○ uses of alternative dispute resolution
○ job search vancouver washington
○ poem of arrival of columbus

Results

```
<topic number="19" type="ambiguous">
    <query>the current</query>
    <description>
        I'm looking for the homepage of The Current, a program
        on Minnesota Public Radio.
    </description>
    <subtopic number="1" type="nav">
        Take me to the homepage of The Current, a program on Minnesota
        Public Radio.
    </subtopic>
    <subtopic number="2" type="nav">
        I'm looking for the homepage of The Current newspaper in New Jersey.
    </subtopic>
    <subtopic number="3" type="nav">
        I want to find the homepage of The Current newspaper in Hartford.
    </subtopic>
    <subtopic number="4" type="nav">
        I want to find the homepage of The Current magazine in San Antonio.
    </subtopic>
</topic>
```

QUERY: job search vancouver washington

DESCRIPTION: I would like to find web page that aggregate job opportunities in the IT industry in Vancouver, Washington.

ty?

ently

eted

em can

ser's intent

empts to

at may be

e of

# Intent-Aware Measures

- Assume there is a probability distribution *P(i|Q)* over intents for a query Q

    – Probability that a randomly-sampled user means intent i when submitting query Q

- The intent-aware version of a measure is its weighted average over this distribution

P(Prado Museum | Q) = 0.35

P(Prado Balboa | Q ) = 0.10

P(Toyota Prado | Q) = 0.45

P(Prado PHP | Q) = 0.08

P(Prado EU | Q) = 0.02

Precision@10-IA = 0.35*0.3 + 0.10*0.2
+ 0.45*0.2 + 0.08*0.1 + 0.02*0.2 = 0.227

# Variance due to Query Intents

- The intents are not fixed
  - may think of them as a sample
    - from some large/infinite universe of possible intents

  - now we have two simultaneous sampling processes
    - what does this (sample x sample) tell us about the (population x population)?

# Mixed Effects Models

- Two sources of variance
  - Queries
  - Intents (within queries)

$$y_{ijk} = \beta_i + b_j + c_{ij} + \varepsilon_{ijk}$$

Value of a metric for system i on query j intent k

Effect of system i

Effect of query j

Effect of sampling intents

Residual Error

$$b_j \sim N(0, \sigma_1^2), \quad c_{ij} \sim N(0, \sigma_2^2), \quad \varepsilon_{ijk} \sim N(0, \sigma^2)$$

# Experimental Design

- Experimental setup
  - How many queries vs. how many intents per query?

| TREC 2010 | Query Effect | Intent Effect |
|-----------|--------------|---------------|
| IA-MAP | 0.0478 | 0.0312 |
| IA-ERR | 0.1429 | 0.0650 |

| TREC 2011 | Query Effect | Intent Effect |
|-----------|--------------|---------------|
| IA-MAP | 0.0707 | 0.0607 |
| IA-ERR | 0.2058 | 0.0973 |

# Conclusions

- Choose your measure carefully

- Choose your experimental setup carefully
  - Put your money where most of the variance comes from
  $\Rightarrow$ Increase statistical power

- Always do significance tests
  - Model all the effects
  - Check your assumptions

- Always take results of tests with a grain of salt
  - Especially when the effect size is low