



Lecture 6: Fat cat

Natural experiments

We're now going to spend a little time on one more example -- This is not a clinical trial but has been termed **a “natural experiment” that leads to an interesting set of questions** we are now prepared to answer

The example comes from a paper that appeared on one of the “journals of record” for statistics -- The author and volume information is given below, and I'll hold off on the full record until we've worked through the problem (the title reveals a bit too much)

The authors, **Ho and Imai have graciously provided us with the data** so you can have a hand at the analysis as well!

Daniel E. Ho and Kosuke Imai,
Journal of the American Statistical Association, Vol. 101, No. 475.

The 2003 recall election

In 2003, California held its **first gubernatorial recall election**, and ultimately the incumbent Democratic Governor Gray Davis was replaced by Republican candidate Arnold Schwarzenegger

To refresh your memories, the California recall system was established in the early 1900s and was introduced to **enact the “will of the people”** who might collectively decide that **an elected official was not doing their job** -- Prior to Davis, many governors had faced the threat of a recall, but in each case had managed to avoid the vote

YOU ARE HERE: LAT Home → Collections → Voting

ADS BY GOOGLE

THE RECALL ELECTION

Gov. Davis Is Recalled; Schwarzenegger Wins

'I Will Not Fail You,' the Republican Victor Promises

October 08, 2003 | Michael Finnegan | Times Staff Writer

Arnold Schwarzenegger won the historic California recall election Tuesday as a tide of voter anger toppled Gray Davis just 11 months after the Democrat had been reelected governor.

In a popular revolt unmatched in the 92 years that Californians have held the power to recall elected officials, voters chose a Republican film star with no government experience to replace an incumbent steeped for three decades in state politics.

Recommend

0



0

Tweet

Submit



Shocking Muscle Pictures
Scientists in Cambridge have discovered a revolutionary new muscle builder. [Read More »](#)

* OK, the ad paired with this particular story is amusing...

2003 recall election

The recall ballot consisted of **two parts** -- In the first you were asked if you thought **Gray Davis should be recalled** and in the second you were asked to **choose a successor** (assuming a majority of voters responded “yes” to the first question)

The list of choices that year, however, was quite long, and in all, **135 candidates qualified to run as a Davis’ replacement** -- Here is a list of their names and how they appeared on sample ballots from assembly districts in Sonoma and Orange Counties

ADAM	ADAMS	ALEXSTJAMES	ANDERSON
ANGELYNE	ARIF	BADIOZAMANI	BAJWA
BEARD	BEYER	BHOLA	BLYCHESTER
BOCK	BRITTON	BROWN	BURTON
BUSTAMANTE	CAMEJO	CARSON	CHAMBERS
CHELI	CLEMENTS	COLEMAN	COOK
CULLENBINE	DAVIS	DOLE	EDWARDS
FARRELL	FEINSTEIN	FLYNT	FONTANES
FORTE	FOSS	FRIEDMAN	GALLAGHER
GORMAN	GOSSE	GREEN	GRISHAM
GRUENER	GUZZARDI	HALL	HAMIDI
HANLON	HENDERSON	HERNANDEZ	HICKEY
HOFFMANN	HUFFINGTON	ISSA	JACKSON
KELLY	KENNEDY	KESSINGER	KIMBALL
KNAPP	KOREVAAR	KUNZMAN	LANE
LEONARD	LEWIS	LOUIE	MACALUSO
MAILANDER	MANNHEIM	MARGOLIN	MARIANO
MARTORANA	MCCARTHY	MCCLAIN	MCCLINTOCK
MCMAHON	MCNEILLY	MEDNICK	MEHR
MILLER	MOBLEY	MOCK	MORTENSEN
MUSILLI	NAVE	NEWMANII	PADILLA
PALMIERI	PAWLIK	PETERS	PINEDA
PRADY	PRICE	QUINN	RAINFORTH
RAMIREZ	RANKEN	RENZ	RICHARDS
RICHTER	RIGHTMYER	ROBINSON	ROSCOE
RUSHFORD	RUSSELL	SAFFORD	SAMS
SCHEIDLE	SCHMIER	SCHWARTZMAN	SCHWARZENEGGER
SIMMONS	SIMON	SMITH	SPRAGUE
SPROUL	STRAUSS	SYLVESTER	TAYLOR
TEMPLIN	TILLEY	TRACY	TSANGARES
UEBERROTH	VALDEZ	VANDEVENTER	VANN
VAUGHN	VO	WALKERC	WALKERM
WALTON	WATTS	WEBER	WEIR
WINTERS	WOZNIAK	ZELLHOEFER	

**Statewide Special Election
Orange County, California**

OFFICIAL BALLOT

October 07, 2003

Instruction Note:

HOW TO VOTE:
To vote, fill in and BLACKEN completely the rectangle in front of any candidate or to the left of the word "YES" or "NO".
Votes for only ONE of the 135 candidates, OR enter a write-in candidate in the space provided.
(Absentee voters should use a dark pen or a #2 pencil.)

Shall GRAY DAVIS be recalled (removed) from the office of Governor?

YES

NO

Candidates to succeed GRAY DAVIS as Governor if he is recalled:
Vote for One

B.E. SMITH
Independent-Lecturer

DAVID RONALD SAMS
Republican-Businessman/Producer/Writer

JAMIE ROSEMARY SAFFORD
Republican-Business Owner

LAWRENCE STEVEN STRAUSS
Democratic-Lawyer/Businessperson/Student

ARNOLD SCHWARZENEGGER
Republican-Actor/Businessman

GEORGE B. SCHWARTZMAN
Independent-Businessman

MIKE SCHMIER
Democratic-Attorney

DARRIN H. SCHEIDLE
Democrat-Businessman/Entrepreneur

BILL SIMON
Republican-Businessman

RICHARD J. SIMMONS
Independent-Attorney/Businessperson

CHRISTOPHER SPROUL
Democrat-Environmental Attorney

RANDALL D. SPRAGUE
Republican-Discrimination Complaint Investigator

TIM SYLVESTER
Democratic-Entrepreneur

JACK LOYD GRISHAM
Independent-Musician/Laborer

JAMES H. GREEN
Democratic-Firefighter Paramedic/Nurse

GARRETT GRIENER
Democratic-High-Tech Entrepreneur

GEROLD LEE GORMAN
Democratic-Engineer

RICH GOSSE
Republican-Educator

LEO GALLAGHER
Independent-Comedian

JOE GIUZZARDI
Democratic-Teacher/Journalist

JON W. ZELLOHOFER
Independent-Businessperson/Consultant

PAUL NAVI
Democratic-Businessman/Entrepreneur

ROBERT C. NEWMAN II
Republican-Psychologist/Farmer

BRIAN TRACY
Independent-Businessperson/Consultant

A. LAVAR TAYLOR
Independent-Attorney

WILLIAM TSANGARES
Republican-Businessperson

PATRICIA G. TILLEY
Independent-Attorney

DIANE BEALL TEMPILIN
American Independent-Attorney/Realtor

MARY "MARY CAREY" COOK
Independent-Adult Film Actress

GARY COLEMAN
Independent-Actor

TODD CARSON
Republican-Rail Estate Developer

PETER MIGUEL CANEJO
Green-Financial Investment Advisor

MICHAEL CHELI
Independent-Businessman

ROBERT CULLENBINE
Democratic-Reified Businessman

D. (LOGAN DARROW) CLEMENTS
Republican-Businessman

S. ISSA
Republican-Engineer

BOB LYNN EDWARDS
Democratic-Attorney

ERIC KOREVAAR
Democratic-Scientist/Businessman

DARRYL L. MOBLEY
Independent-Businessman/Entrepreneur

KELLY P. KIMBALL
Democratic-Business Executive

D.E. KESSINGER
Democratic-Paralegal/Property Manager

EDWARD "ED" KENNEDY
Democratic-Businessman/Educator

TREK THUNDER KELLY
Independent-Business Executive/Artist

JERRY KUNZMAN
Independent-Chief Executive Officer

PETER V. UEBERTH
Republican-Businessman/Olympics Advisor

BILL PRADY
Democratic-Television Writer/Producer

LEONARD PADILLA
Independent-Law School President

DARIN PRICE
Natural Law-University Chemistry Instructor

GREGORY J. PAWLIK
Independent-Nationalist/Businessman

JONATHAN P. PAUL
Independent-Golf Professional

DENNIS DUGGAN MCMAHON
Republican-Banker

ROBERT C. MANNHEIM
Democratic-Medical Doctor

FRANK A. MACALUSO, JR.
Democratic-Physician/Medical Doctor

PAUL MARIANO
Democratic-Attorney

NICKIE MCNEILLY
Independent-Used Car Dealer

MIKE P. MCCARTHY
Independent

BOB MCCLAIN
Independent-Civil Engineer

CHARLES "CHUCK" PINEDA, JR.
Democratic-State Hearing Officer

HEATHER PETERS
Republican-Mediator

ROBERT "BUTCH" DOLE
Republican-Small Business Owner

SCOTT DAVIS
Independent-Business Owner

RONALD J. FRIEDMAN
Independent-Physician

GENE FORTE
Republican-Executive Recruit/Entrepreneur

DIANA FOSS
Democrat

LORRAINE (ABNER ZURD)
Democratic-Film Maker

FONTANES
Independent

WARREN FARRELL
Democratic-Fathers Issues Author

DAN FEINSTEIN
Democrat

LARRY FLYNT
Democrat-Publisher

VAN VO
Republican-Radio Producer/Businessman

JAMES M. VANDEVERTER, JR.
Republican-Saleman/Businessman

PAUL W. VANN
Republican-Financial Planner

BILL VAUGHN
Democrat-Structural Engineer

MARC VALDEZ
Democratic-Air Pollution Scientist

MOHAMMAD ARIF
Independent-Businessman

ANGELYNE
Independent-Entertainer

DOUGLAS ANDERSON
Republican-Mortgage Broker

IRIS ADAM
Natural Law-Business Analyst

BROOKE ADAMS
Independent-Business Executive

ALEX ST. JAMES
Republican-Public Policy Strategist

JIM HOFFMANN
Republican-Teacher

KEN HAMIDI
Libertarian-Healthcare District Director

SARAH ANN HANLON
Independent-Businesswoman

NINA A. HALL
Green-Custom Denture Manufacturer

JOHN J. "JACK" HICKNEY
Libertarian-Healthcare District Director

RALPH A. HERNANDEZ
Democratic-District Attorney Inspector

C. STEPHEN HENDERSON
Independent-Teacher

ARIANNA HUFFINGTON
Independent-Author/Columnist/Mother

ART BROWN
Democrat-Film Writer/Director

JOEL BRITTON
Independent-Reindeer Meat Packer

AUDIE BOCK
Democratic-Education/Social Businesswoman

VIK S. BAJWA
Democratic-Businessman/father/Entrepreneur

BADI BADIOZAMANI
Independent-Entrepreneur/Author/Executive

VIP BHOLA
Republican-Attorney/Businesswoman

JOHN W. BEARD
Independent-Businessman

ED BEYER
Republican-Chief Operations Officer

JOHN CHRISTOPHER BURTON
Independent-Civil Rights Lawyer

CRUZ M. BUSTAMANTE
Democratic-Lesbian/Governor

CHERYL BLY-CHESTER
Republican-Businesswoman/Environmental Engineer

S. ISSA
Independent-Student

LINGEL H. WINTERS
Democratic-Consumer Business Attorney

C.T. WEBER
Independent-Peerless

MAURICE WALKER
Green-Real Estate Appraiser

CHUCK WALKER
Republican-Business Intelligence Analyst

NATHAN WHITE CLOUD WALTON
Independent

MICHAEL WALKER
Green-Real Estate Appraiser

JIM WEIR
Democrat-Community College Teacher

BRYAN QUINN
Republican-Businessman

MICHAEL JACKSON
Republican-State Project Manager

JOHN "JACK" MORTENSEN
Democratic-Contractor/Businessman

ERIC KOREVAAR
Democratic-Scientist/Businessman

S. ISSA
Independent-Engineer

BOB LYNN EDWARDS
Democratic-Attorney

ERIC KOREVAAR
Democratic-Scientist/Businessman

B. T. WEIR
Independent

<input type="checkbox

**Statewide Special Election
Orange County, California
October 07, 2003**

OFFICIAL BALLOT

Instruction Note:

HOW TO VOTE:
To vote, fill in and BLACKEN completely the rectangle to the left of any candidate or to the left of the word "YES" or "NO".

Vote for only ONE of the 135 candidates, OR enter a write-in candidate in the space provided.

Use only the special marking device provided.

(Absentee voters should use a dark pen or a #2 pencil.)

**Shall GRAY DAVIS be recalled (removed)
from the office of Governor?**

YES

NO

**Candidates to succeed GRAY DAVIS as
Governor if he is recalled:**

Vote for One

B.E. SMITH

Independent-Lecturer

DAVID RONALD SAMS

Republican-Businessman/Producer/Writer

JAMIE ROSEMARY SAFFORD

Republican-Business Owner

LAWRENCE STEVEN STRAUSS

Democratic-Lawyer/Businessperson/Student

ARNOLD SCHWARZENEGGER

Republican-Actor/Businessman

GEORGE B. SCHWARTZMAN

Independent-Businessman

MIKE SCHMIER

Democratic-Attorney

DARRIN H. SCHEIDLE

Democratic-Businessman/Entrepreneur

BILL SIMON

Republican-Businessman

RICHARD J. SIMMONS

Independent-Attorney/Businessperson

CHRISTOPHER SPROUL

Democratic-Environmental Attorney

RANDALL D. SPRAGUE

Republican-Discrimination Complaint
Investigator

TIM SYLVESTER

Democratic-Entrepreneur

- | | |
|---|---|
| <input type="checkbox"/> STEPHEN L. KNAPP | Republican-Engineer |
| <input type="checkbox"/> KELLY P. KIMBALL | Democratic-Business Executive |
| <input type="checkbox"/> D.E. KESSINGER | Democratic-Paralegal/Property Manager |
| <input type="checkbox"/> EDWARD "ED" KENNEDY | Democratic-Businessman/Educator |
| <input type="checkbox"/> TREK THUNDER KELLY | Independent-Business Executive/Artist |
| <input type="checkbox"/> JERRY KUNZMAN | Independent-Chief Executive Officer |
| <input type="checkbox"/> PETER V. UEBERROTH | Republican-Businessman/Olympics Advisor |
| <input type="checkbox"/> BILL PRADY | Democratic-Television Writer/Producer |
| <input type="checkbox"/> DARIN PRICE | Natural Law-University Chemistry Instructor |
| <input type="checkbox"/> GREGORY J. PAWLICK | Republican-Realtor/Businessman |
| <input type="checkbox"/> LEONARD PADILLA | Independent-Law School President |
| <input type="checkbox"/> RONALD JASON PALMIERI | Democratic-Gay Rights Attorney |
| <input type="checkbox"/> CHARLES "CHUCK" PINEDA JR. | Democratic-State Hearing Officer |
| <input type="checkbox"/> HEATHER PETERS | Republican-Mediator |
| <input type="checkbox"/> ROBERT "BUTCH" DOLE | Republican-Small Business Owner |
| <input type="checkbox"/> SCOTT DAVIS | Independent-Business Owner |
| <input type="checkbox"/> RONALD J. FRIEDMAN | Independent-Physician |
| <input type="checkbox"/> GENE FORTE | Republican-Executive Recruiter/Entrepreneur |
| <input type="checkbox"/> DIANA FOSS | Democratic- |
| <input type="checkbox"/> LORRAINE (ABNER ZURD) FONTANES | Democratic-Film Maker |
| <input type="checkbox"/> WARREN FARRELL | Democratic-Fathers' Issues Author |
| <input type="checkbox"/> DAN FEINSTEIN | Democratic- |
| <input type="checkbox"/> LARRY FLYNT | Democratic-Publisher |

- | | |
|--|--|
| <input type="checkbox"/> DARRYL L. MOBLEY | Independent-Businessman/Entrepreneur |
| <input type="checkbox"/> JEFFREY L. MOCK | Republican-Business Owner |
| <input type="checkbox"/> BRUCE MARGOLIN | Democratic-Marijuana Legalization Attorney |
| <input type="checkbox"/> GINO MARTORANA | Republican-Restaurant Owner |
| <input type="checkbox"/> PAUL MARIANO | Democratic-Attorney |
| <input type="checkbox"/> ROBERT C. MANNHEIM | Democratic-Retired Businessperson |
| <input type="checkbox"/> FRANK A. MACALUSO, JR. | Democratic-Physician/Medical Doctor |
| <input type="checkbox"/> PAUL "CHIP" MAJLANDER | Democratic-Golf Professional |
| <input type="checkbox"/> DENNIS DUGGAN MCMAHON | Republican-Banker |
| <input type="checkbox"/> MIKE MCNEILLY | Republican-Artist |
| <input type="checkbox"/> MIKE P. MCCARTHY | Independent-Used Car Dealer |
| <input type="checkbox"/> BOB MCCLAIN | Independent-Civil Engineer |
| <input type="checkbox"/> TOM MCCLINTOCK | Republican-State Senator |
| <input type="checkbox"/> JONATHAN MILLER | Democratic-Small Business Owner |
| <input type="checkbox"/> CARL A. MEHR | Republican-Businessman |
| <input type="checkbox"/> SCOTT A. MEDNICK | Democratic-Business Executive |
| <input type="checkbox"/> DORENE MUSILLI | Republican-Parent/Educator/Businesswoman |
| <input type="checkbox"/> VAN VO | Republican-Radio Producer/Businessman |
| <input type="checkbox"/> PAUL W. VANN | Republican-Financial Planner |
| <input type="checkbox"/> JAMES M. VANDEVENTER, JR. | Republican-Salesman/Businessman |
| <input type="checkbox"/> BILL VAUGHN | Democratic-Structural Engineer |
| <input type="checkbox"/> MARC VALDEZ | Democratic-Air Pollution Scientist |
| <input type="checkbox"/> MOHAMMAD ARIF | Independent-Businessman |

MEASURES SUBMITTED TO THE VOTERS

STATE

Proposition 53

**FUNDS DEDICATED FOR STATE AND
LOCAL INFRASTRUCTURE.
LEGISLATIVE CONSTITUTIONAL
AMENDMENT.**

Generally dedicates up to 3% of General Fund revenues annually to fund state and local (excluding school and community college) infrastructure projects. Fiscal Impact: Dedication of General Fund revenues for state and local infrastructure. Potential transfers of \$850 million in 2006-07, increasing to several billions of dollars in future years, under specified conditions.

YES

NO

Proposition 54

**CLASSIFICATION BY RACE, ETHNICITY,
COLOR, OR NATIONAL ORIGIN.
INITIATIVE CONSTITUTIONAL
AMENDMENT.**

Prohibits state and local governments from classifying any person by race, ethnicity, color, or national origin. Various exemptions apply. Fiscal Impact: The measure would not result in a significant fiscal impact on state and local governments.

YES

NO

SAMPLE BALLOT

OFFICIAL BALLOT

Statewide Special Election

Sonoma County

October 7, 2003

This ballot stub shall be removed and retained by the voter.

**MARK YOUR CHOICE(S)
IN THIS MANNER ONLY:** 

I HAVE VOTED—HAVE YOU?

**MARK YOUR CHOICE(S)
IN THIS MANNER ONLY:** 

STATE		
Shall GRAY DAVIS be recalled (removed) from the office of Governor?	Yes	No
Candidates to succeed GRAY DAVIS as Governor if he is recalled.	Vote for One	
KURT E. "TACHIKAZE" RIGHTEMYER, Independent Middleweight Sumo Wrestler		
DANIEL W. RICHARDS, Republican Businessman	ROBERT C. MANNHEIM, Democratic Retired Businessperson	FRANK A. MACALUSO, JR., Democratic Physician/Medical Doctor
KEVIN RICHTER, Republican Information Technology Manager	PAUL "CHIP" MAILANDER, Democratic Golf Professional	DENNIS DUGGAN McMAHON, Republican Banker
REVA RENEE RENZ, Republican Small Business Owner	MIKE MCNEILLY, Republican Artist	MIKE P. McCARTHY, Independent Used Car Dealer
SHARON RUSHFORD, Independent Businesswoman	BOB MCCLAIN, Independent Civil Engineer	TOM MCCINTOCK, Republican State Senator
GEORGY RUSSELL, Democratic Software Engineer	JONATHAN MILLER, Democratic Small Business Owner	CARL A. MEHR, Republican Businessman
MICHAEL J. WOZNIAK, Democratic Retired Police Officer	SCOTT A. MEDNICK, Democratic Business Executive	DORENE MUSSILLI, Republican Parent/Educator/Businesswoman
DANIEL WATTS, Green College Student	VAN VO, Republican Radio Producer/Businessman	PAUL W. VANN, Republican Financial Planner
NATHAN WHITECLOUD WALTON, Independent Student	JAMES M. VANDEVENTER, JR., Republican Salesman/Businessman	BILL VAUGHN, Democratic Structural Engineer
MAURICE WALKER, Green Real Estate Appraiser	MARC VALDEZ, Democratic Air Pollution Scientist	MOHAMMAD ARIF, Independent Businessman
CHUCK WALKER, Republican Business Intelligence Analyst	ANGELYNE, Independent Entertainer	DOUGLAS ANDERSON, Republican Mortgage Broker
LINGEL H. WINTERS, Democratic Consumer Business Attorney	JIM WEIR, Democratic Community College Teacher	IRIS ADAM, Natural Law Business Analyst
C.T. WEBER, Peace and Freedom Labor Official/Analyst	BRYAN QUINN, Republican Businessman	BROOKE ADAMS, Independent Business Executive
JIM WEIR, Democratic Community College Teacher	MICHAEL JACKSON, Republican Satellite Project Manager	ALEX-ST. JAMES, Republican Public Policy Strategist
BRYAN QUINN, Republican Businessman	JOHN "JACK" MORTENSEN, Democratic Contractor/Businessman	JIM HOFFMANN, Republican Teacher
DARRYL L. MOBLEY, Independent Businessman/Entrepreneur	DARRYL L. MOBLEY, Independent Businessman	KEN HAMIDI, Libertarian State Tax Officer
JEFFREY L. MOCK, Republican Business Owner	BRUCE MARGOLIN, Democratic Marijuana Legalization Attorney	
GINO MARTORANA, Republican Restaurant Owner	GINO MARTORANA, Republican Restaurant Owner	
PAUL MARIANO, Democratic Attorney	PAUL MARIANO, Democratic Attorney	

49-A007R **CONTINUED OTHER SIDE**

(CANDIDATES CONTINUED)		
ROBERT C. MANNHEIM, Democratic Retired Businessperson	FRANK A. MACALUSO, JR., Democratic Physician/Medical Doctor	PAUL "CHIP" MAILANDER, Democratic Golf Professional
DENNIS DUGGAN McMAHON, Republican Banker	MIKE MCNEILLY, Republican Artist	MIKE P. McCARTHY, Independent Used Car Dealer
MIKE P. McCARTHY, Independent Used Car Dealer	BOB MCCLAIN, Independent Civil Engineer	TOM MCCINTOCK, Republican State Senator
MIKE MCNEILLY, Republican Artist	JONATHAN MILLER, Democratic Small Business Owner	CARL A. MEHR, Republican Businessman
MIKE P. McCARTHY, Independent Used Car Dealer	SCOTT A. MEDNICK, Democratic Business Executive	DORENE MUSSILLI, Republican Parent/Educator/Businesswoman
BOB MCCLAIN, Independent Civil Engineer	VAN VO, Republican Radio Producer/Businessman	PAUL W. VANN, Republican Financial Planner
TOM MCCINTOCK, Republican State Senator	JAMES M. VANDEVENTER, JR., Republican Salesman/Businessman	BILL VAUGHN, Democratic Structural Engineer
JONATHAN MILLER, Democratic Small Business Owner	MARC VALDEZ, Democratic Air Pollution Scientist	MOHAMMAD ARIF, Independent Businessman
CARL A. MEHR, Republican Businessman	ANGELYNE, Independent Entertainer	DOUGLAS ANDERSON, Republican Mortgage Broker
SCOTT A. MEDNICK, Democratic Business Executive	IRIS ADAM, Natural Law Business Analyst	
DORENE MUSSILLI, Republican Parent/Educator/Businesswoman	BROOKE ADAMS, Independent Business Executive	
PAUL W. VANN, Republican Financial Planner	ALEX-ST. JAMES, Republican Public Policy Strategist	
BILL VAUGHN, Democratic Structural Engineer	JIM HOFFMANN, Republican Teacher	
MOHAMMAD ARIF, Independent Businessman	KEN HAMIDI, Libertarian State Tax Officer	

49-A008R **CONTINUED NEXT CARD**

A

A

A

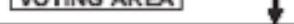
Sonoma County

October 7, 2003

This ballot stub shall be removed and retained by the voter.

MARK YOUR CHOICE(S)
IN THIS MANNER ONLY:

VOTING AREA



STATE

Shall GRAY DAVIS be recalled (removed) from the office of Governor?

Yes

No

Candidates to succeed GRAY DAVIS as Governor if he is recalled.

Vote for One

KURT E. "TACHIKAZE" RIGHTMYER, Independent
Middleweight Sumo Wrestler

DANIEL W. RICHARDS, Republican
Businessman

KEVIN RICHTER, Republican
Information Technology Manager

REVA RENEE RENZ, Republican
Small Business Owner

SHARON RUSHFORD, Independent
Businesswoman

GEORGY RUSSELL, Democratic
Software Engineer

MICHAEL J. WOZNIAK, Democratic
Retired Police Officer

DANIEL WATTS, Green
College Student

NATHAN WHITECLOUD WALTON, Independent
Student

MAURICE WALKER, Green
Real Estate Appraiser

CHUCK WALKER, Republican
Business Intelligence Analyst

LINGEL H. WINTERS, Democratic
Consumer Business Attorney

I HAVE VOTED—HAVE YOU?

MARK YOUR CHOICE(S)
IN THIS MANNER ONLY:

VOTING AREA



(CANDIDATES CONTINUED)

ROBERT C. MANNHEIM, Democratic
Retired Businessperson

FRANK A. MACALUSO, JR., Democratic
Physician/Medical Doctor

PAUL "CHIP" MAILANDER, Democratic
Golf Professional

DENNIS DUGGAN MCMAHON, Republican
Banker

MIKE MCNEILLY, Republican
Artist

MIKE P. MCCARTHY, Independent
Used Car Dealer

BOB MCCLAIN, Independent
Civil Engineer

TOM MCCLINTOCK, Republican
State Senator

JONATHAN MILLER, Democratic
Small Business Owner

CARL A. MEHR, Republican
Businessman

SCOTT A. MEDNICK, Democratic
Business Executive

DORENE MUSILLI, Republican
Parent/Educator/Businesswoman

VAN VO, Republican
Radio Producer/Businessman

PAUL W. VANN, Republican
Financial Planner

JAMES M. VANDEVENTER, JR., Republican
Consumer Business Attorney

San

2003 recall election

There were so many names that most voting districts had to **split them up into multiple pages** -- In addition to that formatting question, what else do you notice about the official ballots on the previous pages?

In particular, what can you tell me about the order of the names?

KURT E. "TACHIKAZE" RIGHTMYER, Independent
Middleweight Sumo Wrestler

DANIEL W. RICHARDS, Republican
Businessman

KEVIN RICHTER, Republican
Information Technology Manager

REVA RENEE RENZ, Republican
Small Business Owner

SHARON RUSHFORD, Independent
Businesswoman

GEORGY RUSSELL, Democratic
Software Engineer

MICHAEL J. WOZNIAK, Democratic
Retired Police Officer

DANIEL WATTS, Green
College Student

NATHAN WHITECLOUD WALTON, Independent
Student

MAURICE WALKER, Green
Real Estate Appraiser

CHUCK WALKER, Republican
Business Intelligence Analyst

LINGEL H. WINTERS, Democratic
Consumer Business Attorney

C.T. WEBER, Peace and Freedom
Labor Official/Analyst

JIM WEIR, Democratic
Community College Teacher

BRYAN QUINN, Republican
Businessman

MICHAEL JACKSON, Republican
Satellite Project Manager

JOHN "JACK" MORTENSEN, Democratic
Contractor/Businessman

DARRYL L. MOBLEY, Independent
Businessman/Entrepreneur

JEFFREY L. MOCK, Republican
Business Owner

BRUCE MARGOLIN, Democratic
Marijuana Legalization Attorney

GINO MARTORANA, Republican
Restaurant Owner

PAUL MARIANO, Democratic
Attorney

ROBERT C. MANNHEIM, Democratic
Retired Businessperson

FRANK A. MACALUSO, JR., Democratic
Physician/Medical Doctor

PAUL "CHIP" MAILANDER, Democratic
Golf Professional

DENNIS DUGGAN MCMAHON, Republican
Banker

MIKE MCNEILLY, Republican
Artist

MIKE P. MCCARTHY, Independent
Used Car Dealer

BOB MCCLAIN, Independent
Civil Engineer

TOM MCCLINTOCK, Republican
State Senator

JONATHAN MILLER, Democratic
Small Business Owner

CARL A. MEHR, Republican
Businessman

SCOTT A. MEDNICK, Democratic
Business Executive

DORENE MUSILLU, Republican
Parent/Educator/Businesswoman

VAN VO, Republican
Radio Producer/Businessman

PAUL W. VANN, Republican
Financial Planner

California alphabet soup

Prior to 1975, the state tended to place incumbents at the top of its ballots -- In that year, however, **the state Supreme Court decided that listing candidates according to incumbency status or even in a straight alphabetical order was unconstitutional**

... the superior court's finding that placement in a top ballot position affords a candidate a substantial advantage over lower-placed candidates is supported by abundant expert testimony introduced at trial and is consistent with parallel findings rendered in similar litigation throughout the country. In light of this finding, we explain that **any procedure which allocates such advantageous positions to a particular class of candidates inevitably discriminates against voters supporting all other candidates**, and accordingly can only be sustained if necessary to further a compelling governmental interest. Applying this test, we conclude that the city [Santa Monica] has demonstrated no compelling interest which necessitates the provision's discriminatory classification scheme and thus we uphold the trial court's determination of invalidity. Finally, with respect to a subsidiary matter, we conclude that the allocation of advantageous ballot positions on the basis of "alphabetical order" is similarly unconstitutional.

California alphabet soup

To address this, the state legislature developed a randomization process in which the Secretary of State creates **a new alphabetical ordering** to be used when sorting candidates

Quite literally, this means putting all **26 letters in a hat (well, something hat-like) and drawing them out one at a time** -- The order in which each letter is drawn specifies its precedence when sorting candidate names

Here is Section 13112 of the California Elections Code...

The Secretary of State shall conduct a drawing of the letters of the alphabet, the result of which shall be known as a randomized alphabet. The procedure shall be as follows:

(a) **Each letter of the alphabet shall be written on a separate slip of paper, each of which shall be folded and inserted into a capsule.** Each capsule shall be opaque and of uniform weight, color, size, shape, and texture. **The capsules shall be placed in a container, which shall be shaken vigorously in order to mix the capsules thoroughly. The container then shall be opened and the capsules removed at random one at a time.** As each is removed, it shall be opened and the letter on the slip of paper read aloud and written down. **The resulting random order of letters constitutes the randomized alphabet, which is to be used in the same manner as the conventional alphabet in determining the order of all candidates in all elections.** For example, if two candidates with the surnames Campbell and Carlson are running for the same office, their order on the ballot will depend on the order in which the letters M and R were drawn in the randomized alphabet drawing.

(b) (1) There shall be six drawings, three in each even-numbered year and three in each odd-numbered year. Each drawing shall be held at 11 a.m. on the date specified in this subdivision. The results of each drawing shall be mailed immediately to each county elections official responsible for conducting an election to which the drawing is applicable, who shall use it in determining the order on the ballot of the names of the candidates for office.

(A) The first drawing under this subdivision shall take place on the 82nd day before the April general law city elections of an even-numbered year, and shall apply to those elections and any other elections held at the same time.

(B) The second drawing under this subdivision shall take place on the 82nd day before the direct primary of an even-numbered year, and shall apply to all candidates on the ballot in that election.

(C) (i) The third drawing under this subdivision shall take place on the 82nd day before the November general election of an even-numbered year, and shall apply to all candidates on the ballot in the November general election.

(ii) In the case of the primary election and the November general election, the Secretary of State shall certify and transmit to each county elections official the order in which the names of federal and state candidates, with the exception of candidates for State Senate and Assembly, shall appear on the ballot. The elections official shall determine the order on the ballot of all other candidates using the appropriate randomized alphabet for that purpose.

(D) The fourth drawing under this subdivision shall take place on the 82nd day before the March general law city elections of each odd-numbered year, and shall apply to those elections and any other elections held at the same time.

(E) The fifth drawing under this subdivision shall take place on the 82nd day before the first Tuesday after the first Monday in June of each odd-numbered year, and shall apply to all candidates on the ballot in the elections held on that date.

(F) The sixth drawing under this subdivision shall take place on the 82nd day before the first Tuesday after the first Monday in November of the odd-numbered year, and shall apply to all candidates on the ballot in the elections held on that date.

(2) In the event there is to be an election of candidates to a special district, school district, charter city, or other local government body at the same time as one of the five major election dates specified in subparagraphs (A) to (F), inclusive, and the last possible day to file nomination papers for the local election would occur after the date of the drawing for the major election date, the procedure set forth in Section 13113 shall apply.

(c) Each randomized alphabet drawing shall be open to the public. At least 10 days prior to a drawing, the Secretary of State shall notify the news media and other interested parties of the date, time, and place of the drawing. The president of each statewide association of local officials with responsibilities for conducting elections shall be invited by the Secretary of State to attend each drawing or send a representative. The state chairman of each qualified political party shall be invited to attend or send a representative in the case of drawings held to determine the order of candidates on the primary election ballot, the November general election ballot, or a special election ballot as provided for in subdivision (d).

California alphabet soup

As an example, suppose the lottery produced this sequence of letters

O G F W N T Q B M U Y R S D K H V E X J Z I P A C L

Then, using this sequence, we can produce the following sorted list of names

GORMAN	GOSSE	GUZZARDI	GRUENER
GREEN	GRISHAM	GALLAGHER	FONTANES
FORTE	FOSS	FRIEDMAN	FEINSTEIN
FARRELL	FLYNT	WOZNIAK	WEBER
WEIR	WINTERS	WATTS	WALTON
WALKERM	WALKERC	NEWMANII	NAVE
TRACY	TSANGARES	TEMPLIN	TILLEY
TAYLOR	QUINN	BOCK	BURTON
BUSTAMANTE	BROWN	BRITTON	BHOLA
BEYER	BEARD	BADIOZAMANI	BAJWA
BLYCHESTER	MOBLEY	MORTENSEN	MOCK
MUSILLI	MEDNICK	MEHR	MILLER
MANNHEIM	MARGOLIN	MARTORANA	MARIANO
MAILANDER	MACALUSO	MCNEILLY	MCMAHON
MCCARTHY	MCCLINTOCK	MCCLAIN	UEBERROTH

California alphabet soup

Or, suppose the lottery produced this sequence of letters

F R I A Z U M Y J N X L V O B D S T E Q K W H C P G

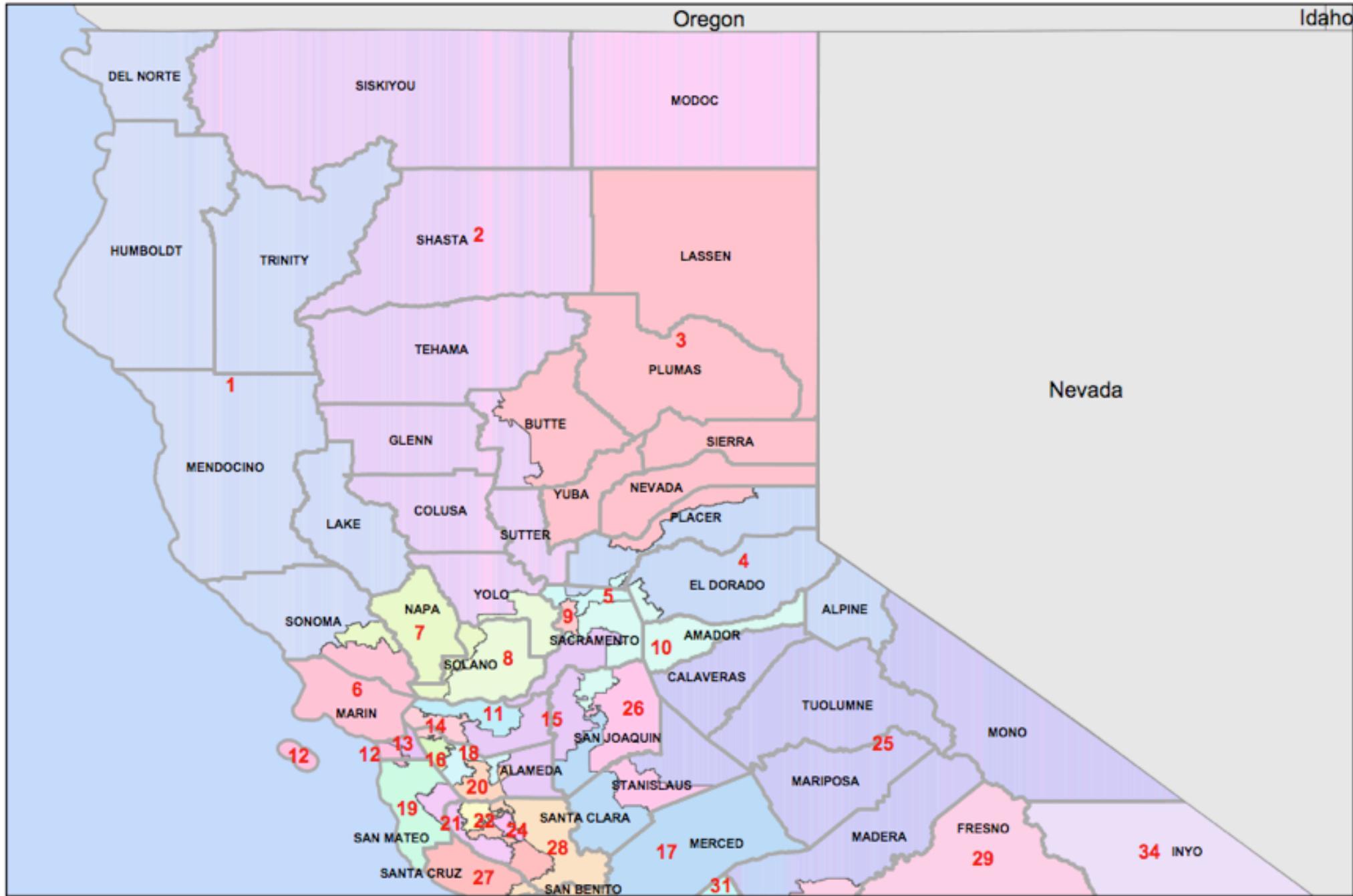
Then, using this sequence, we can produce the following sorted list of names

FRIEDMAN	FARRELL	FLYNT	FORTE
FONTANES	FOSS	FEINSTEIN	RICHARDS
RICHTER	RIGHTMYER	RAINFORTH	RAMIREZ
RANKEN	RUSSELL	RUSHFORD	ROBINSON
ROSCOE	RENZ	ISSA	ARIF
ANDERSON	ANGELYNE	ALEXSTJAMES	ADAM
ADAMS	ZELLHOFER	UEBERROTH	MILLER
MARIANO	MARTORANA	MARGOLIN	MAILANDER
MANNHEIM	MACALUSO	MUSILLI	MORTENSEN
MOBLEY	MOCK	MEDNICK	MEHR
MCMAHON	MCNEILLY	MCCARTHY	MCCLINTOCK
MCCLAIN	JACKSON	NAVE	NEWMANII
LANE	LOUIE	LEONARD	LEWIS
VAUGHN	VANN	VANDEVENTER	VALDEZ
VO	BRITTON	BROWN	BAJWA

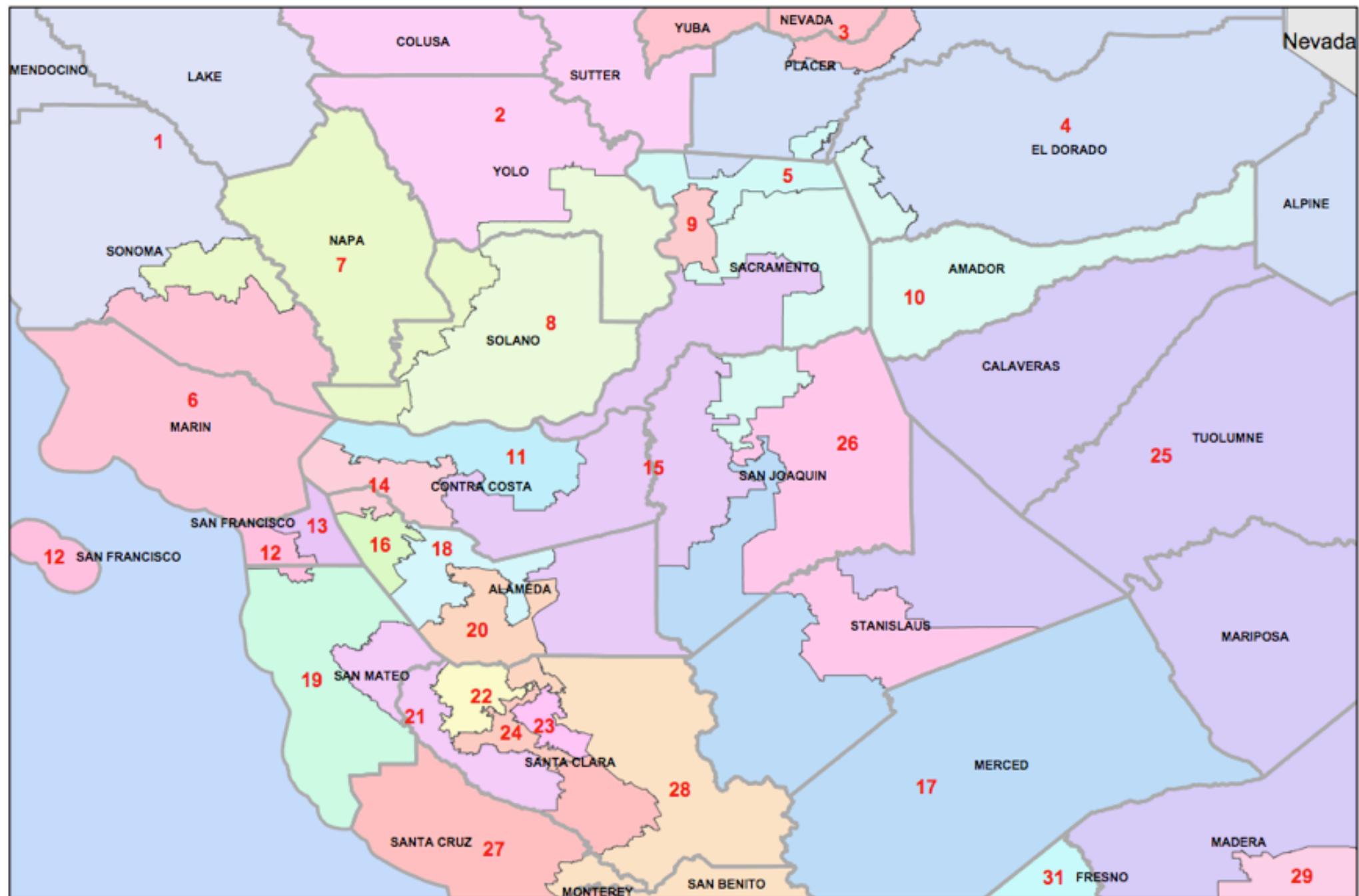
Randomize and rotate

Once an alphabet is set, the candidates ordered but only the first assembly district uses this as its ballot order -- For each subsequent assembly district (we have 80 districts in the 2003 election), **the name at the top of the list is moved to the bottom and all the other names are moved up one place**

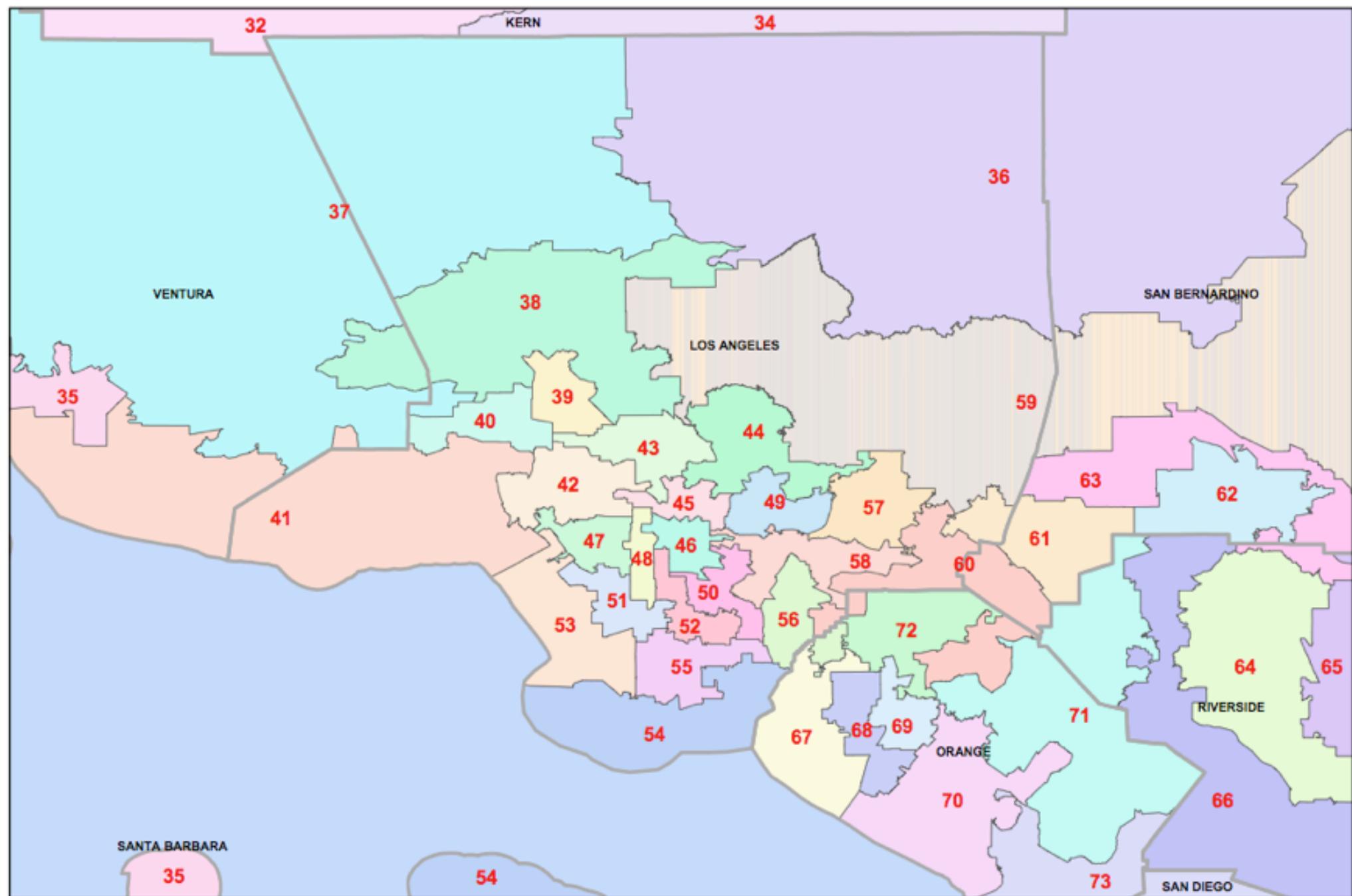
The idea behind this rotation is that candidates each spent time near the top of a ballot in at least some of the districts -- This all seems sensible when **the number of candidates is small**



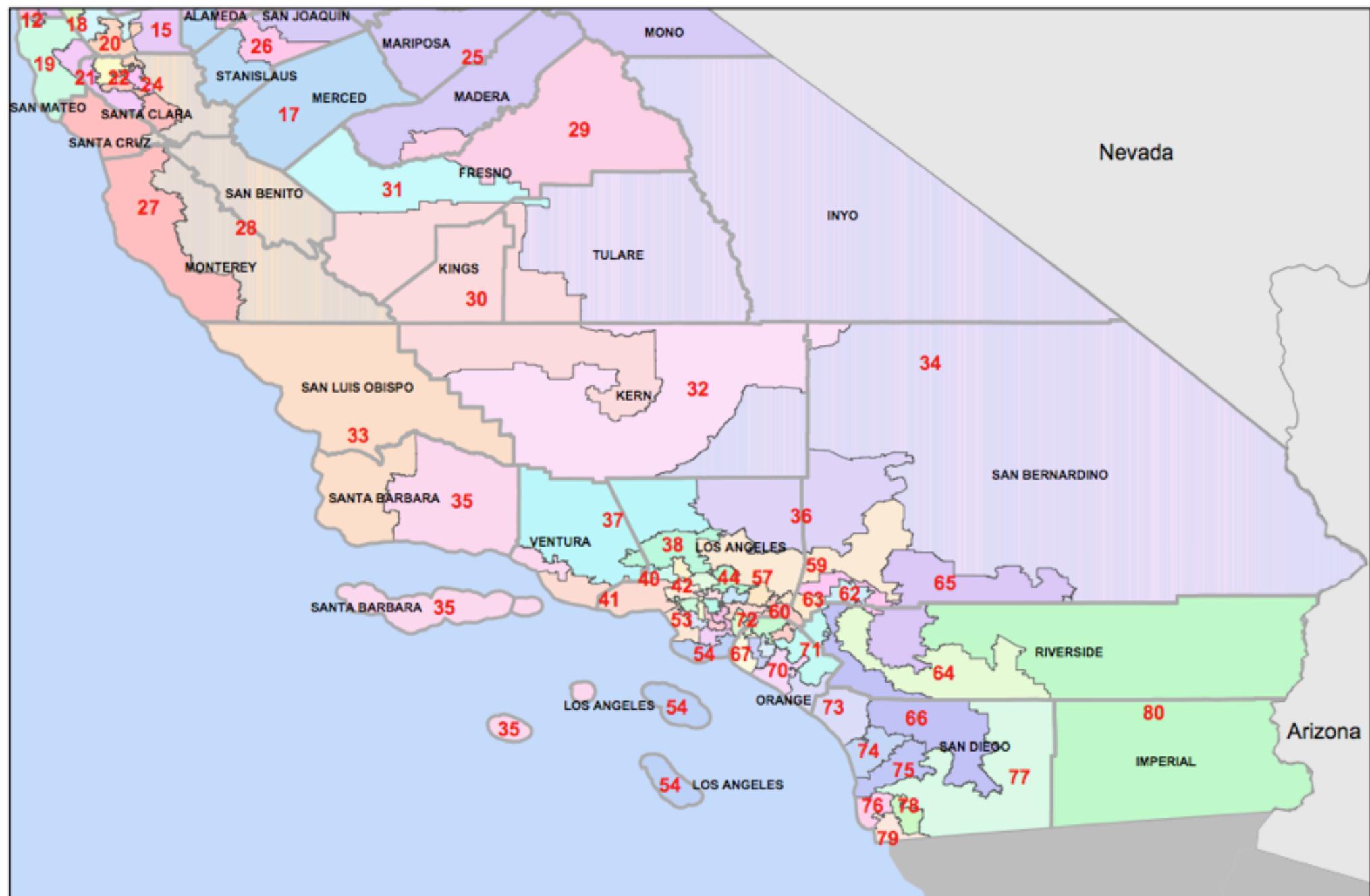
© Copyright 2002, California Voter Foundation, www.calvoter.org, All Rights Reserved.



© Copyright 2002, California Voter Foundation, www.calvoter.org, All Rights Reserved.



© Copyright 2002, California Voter Foundation, www.calvoter.org. All Rights Reserved.



Randomize and rotate

Assembly districts intersect counties -- Some districts contain multiple counties (District 1 contains Del Norte, Humboldt, Lake, Mendocino, Sonoma and Trinity counties) while some counties are split among many districts (Alameda county is in Districts 14, 15, 16, 18, and 20)

While the order of candidates' names is set per assembly district, **the ballot format (the number of names per page) is set by each county** -- In all there are 158 district-county combinations

2003 recall election

The alphabet used in 2003 was

R W Q O J M V A H B S G Z X N T C I E K U P D Y F L

The order and the breakdown by page number for the first district (in Del Norte County) is given on the next slide and then the breakdown by page number for the last district (in Riverside County) is given after that

District 1,
Del Norte county

Page 1	ROBINSON RAINFORTH RENZ WATTS WINTERS JACKSON	ROSCOE RIGHTMYER RUSHFORD WALTON WEBER MORTENSEN	RAMIREZ RICHARDS RUSSELL WALKERM WEIR MOBLEY	RANKEN RICHTER WOZNIAK WALKERC QUINN
Page 2	MOCK MANNHEIM MCNEILLY MILLER VO VALDEZ	MARGOLIN MACALUSO MCCARTHY MEHR VANN ARIF	MARTORANA MAILANDER MCCLAIN MEDNICK VANDEVENTER ANGELYNE	MARIANO MCMAHON MCCLINTOCK MUSILLI VAUGHN ANDERSON
Page 3	ADAM HAMIDI HERNANDEZ BRITTON BHOLA BUSTAMANTE	ADAMS HANLON HENDERSON BOCK BEARD BLYCHESTER	ALEXSTJAMES HALL HUFFINGTON BAJWA BEYER SMITH	HOFFMANN HICKEY BROWN BADIOZAMANI BURTON SAMS
Page 4	SAFFORD SCHMIER SPROUL GREEN GALLAGHER NEWMANII	STRAUSS SCHEIDLE SPRAGUE GRUENER GUZZARDI TRACY	SCHWARZENEGGER SIMON SYLVESTER GORMAN ZELLHOFER TAYLOR	SCHWARTZMAN SIMMONS GRISHAM GOSSE NAVE TSANGARES
Page 5	TILLEY CARSON CULLENBINE KOREVAAR KENNEDY PRADY	TEMPLIN CAMEJO CLEMENTS KNAPP KELLY PRICE	COOK CHAMBERS ISSA KIMBALL KUNZMAN PAWLIK	COLEMAN CHELI EDWARDS KESSINGER UEBERROTH PADILLA
Page 6	PALMIERI DAVIS FONTANES LOUIE	PINEDA FRIEDMAN FARRELL LANE	PETERS FORTE FEINSTEIN LEWIS	DOLE FOSS FLYNT LEONARD

District 80,
Riverside County

Page 1	SPROUL GREEN GALLAGHER NEWMANII TILLEY	SPRAGUE GRUENER GUZZARDI TRACY TEMPLIN	SYLVESTER GORMAN ZELLHOEFER TAYLOR COOK	GRISHAM GOSSE NAVE TSANGARES COLEMAN
	CARSON CULLENBINE KOREVAAR KENNEDY PRADY	CAMEJO CLEMENTS KNAPP KELLY PRICE	CHAMBERS ISSA KIMBALL KUNZMAN PAWLIK	CHELI EDWARDS KESSINGER UEBERROTH PADILLA
	PALMIERI DAVIS FONTANES LOUIE ROBINSON	PINEDA FRIEDMAN FARRELL LANE ROSCOE	PETERS FORTE FEINSTEIN LEWIS RAMIREZ	DOLE FOSS FLYNT LEONARD RANKEN
	RAINFORTH RENZ WATTS WINTERS JACKSON	RIGHTMYER RUSHFORD WALTON WEBER MORTENSEN	RICHARDS RUSSELL WALKERM WEIR MOBLEY	RICHTER WOZNIAK WALKERC QUINN MOCK
	MARGOLIN MACALUSO MCCARTHY MEHR VANN	MARTORANA MAILANDER MCCLAIN MEDNICK VANDEVENTER	MARIANO MCMAHON MCCLINTOCK MUSILLI VAUGHN	MANNHEIM MCNEILLY MILLER VO VALDEZ
	ARIF ADAMS HANLON HENDERSON BOCK	ANGELYNE ALEXSTJAMES HALL HUFFINGTON BAJWA	ANDERSON HOFFMANN HICKEY BROWN BADIOZAMANI	ADAM HAMIDI HERNANDEZ BRITTON BHOLA
	BEARD BLYCHESTER STRAUSS SCHEIDLE	BEYER SMITH SCHWARZENEGGER SIMON	BURTON SAMS SCHWARTZMAN SIMMONS	BUSTAMANTE SAFFORD SCHMIER

Location, location, location

With so many candidates on the ballot, the division into **multiple pages** was inevitable -- In all, **121 of the 158 district-county pairs produced ballots with more than one page**

Recalling the Supreme Court's decision that ballot placement has an effect on voting behavior, the rotation process was meant to **even out the number of times people are banished to the back pages**

We can have a look at how often each of the 135 candidates were on the first page of a ballot in the 121 district-county combinations...

```
# read in a csv (comma separated values) file provided to us by prof daniel ho
# at stanford university

> ballot <- read.csv(url("http://www.stat.ucla.edu/~cocteau/stat105/data/ballot.csv"))

> table(ballot$Schwarzenegger)

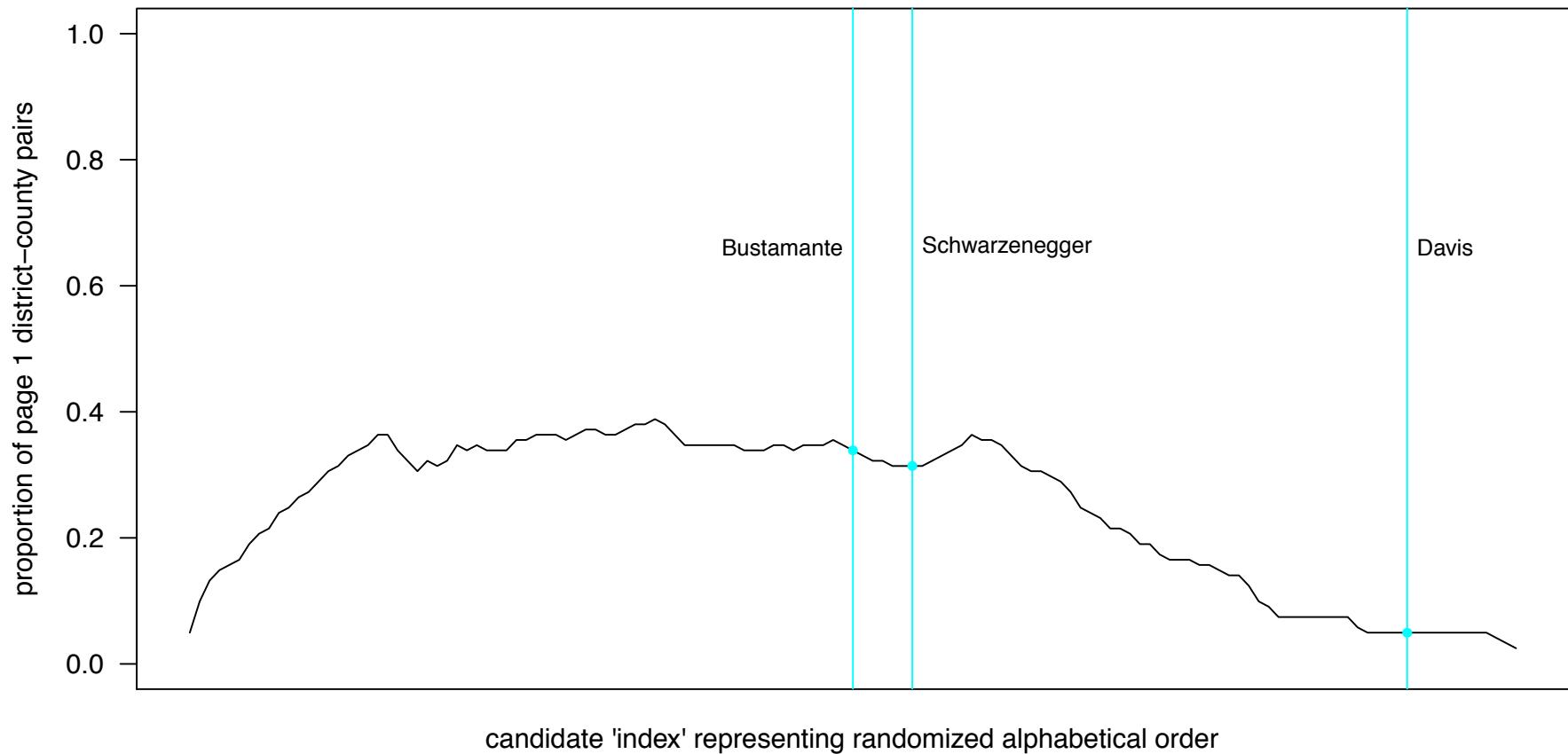
 1   2   3   4   6   7
38  29  30  16   2   6

> table(ballot$Bustamante)

 1   2   3   4   6   7
41  31  35   3   2   9

> table(ballot$Davis)

 1   2   3   4   5   6   7  11
 6  24  11  22  33  21   3   1
```

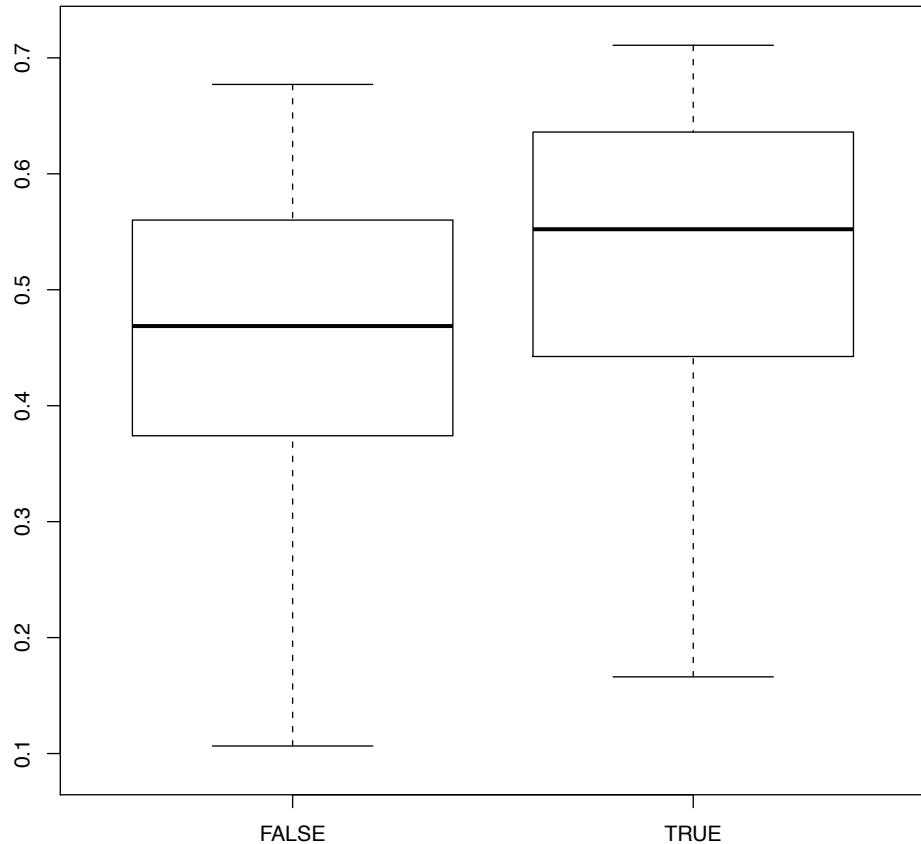


Location, location, location

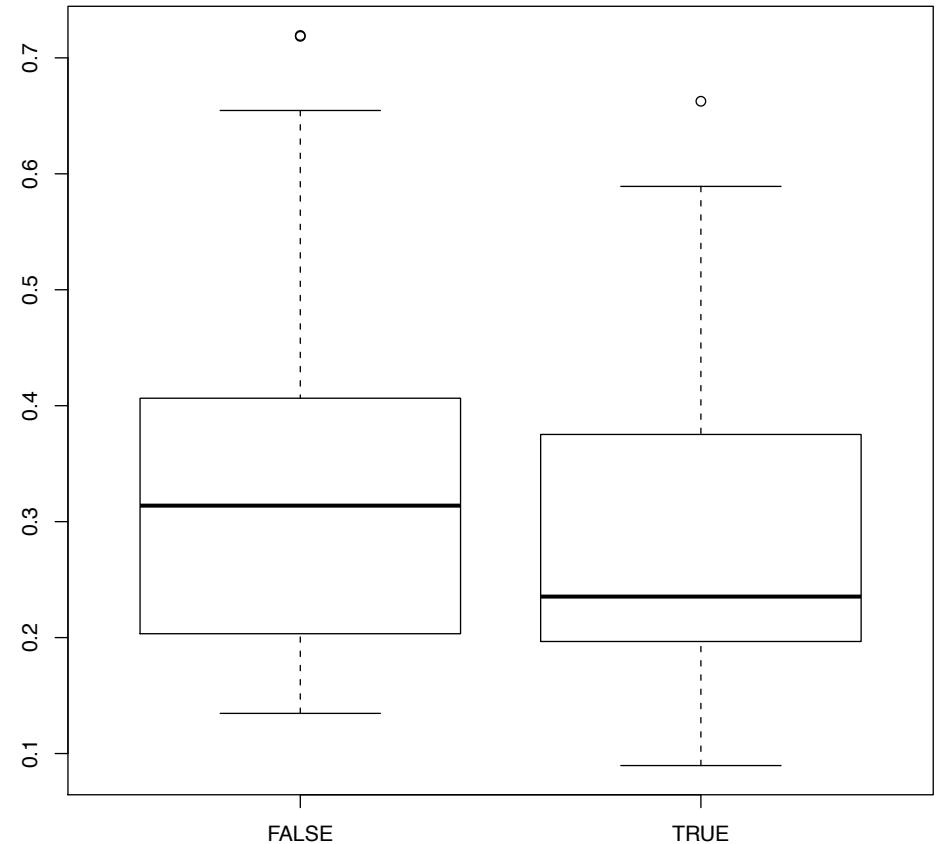
Given the number of candidates in 2003, each was on the front page in less than half of the county-district pairs -- We can reasonably ask whether or not **being on the first page of a ballot has an effect** on the share of the vote a candidate receives

How might we measure this?

Shares for Schwarzenegger against page 1 status



Shares for Bustamante against page 1 status



```
> shares <- read.csv(url("http://www.stat.ucla.edu/~cocteau/stat105/data/shares.csv"))

> head(shares$Schwarzenegger)
[1] 0.5497971 0.4709309 0.3719113 0.3448414 0.6331666 0.6420732

> head(ballot$Schwarzenegger)
[1] 4 4 4 4 4 1

> boxplot(shares$Schwarzenegger ~ (ballot$Schwarzenegger == 1))
```

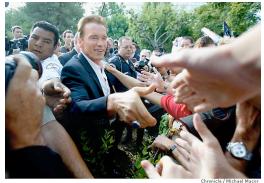
Vote shares

One natural measure is just the difference between the average share (proportion) of the vote a candidate received in districts where they were on the first page of the ballot and the average computed over districts where they were on a later page

In symbols, if we let y_i denote one candidate's vote share in district i , then we can capture the effect of being on the first page with

$$T = (\text{mean } y_i \text{ for districts } i \text{ where candidate is listed on page 1}) - (\text{mean } y_i \text{ for districts } i \text{ where candidate is on later page})$$

Vote shares: Examples



As an example, **Arnold Schwarzenegger** received an average of 52% of the vote in districts where he was listed on the first page of the ballot, but only 46% when he was on page 2 or higher -- That is, $T = 0.05$



Cruz Bustamante received an average of 29% of the vote in districts where he was listed on the first page of the ballot, and an average of 33% of the vote in the remaining districts, so that $T = -0.04$



Ariana Huffington (who eventually dropped out of the race although her name remained on the ballot) received an average of 0.54% of the vote in districts where she was on page one of the ballot, but an average of 0.52% in other districts, implying $T = 0.0002$



Finally, **Bruce Margolin** was a relatively small player earning an average of 0.2% of the vote when on the first page, but 0.1% when listed on later pages so that $T = 0.001$

Starting to look familiar

What do we make of these differences? Does our analysis end with observing the differences or can we say more? Given that we've spent a week talking about "statistically significant" differences, is there some way for us apply that framework here?

Hint: I wouldn't have burned 30 slides on this story if we couldn't do more! So, what do you think?

Re-randomization

Given that the State of California creates **a randomized alphabet** to start the ballot placement process, we have a kind of “**natural experiment**” that we can exploit to assess whether the differences T we see for a given candidate could be **purely the result of randomization**

Let’s start with the **null hypothesis that page placement has no effect on the vote share earned by a candidate** -- Under this assumption, we can **generate new alphabets, create new ballot layouts and compute a new difference in mean vote shares**

Repeating this many times for a given candidate, we create a null distribution for the difference in mean vote shares T , to which we can **compare the difference that was actually recorded on election day...**

One iteration

In the original alphabet

RWQOJMVAHBSGXNTCIEKUPDYFL

Schwarzenegger appeared in the 74th position in District 1 -- After rotation, here are the page numbers he appeared on in the 121 districts with multi-page ballots

His average vote share in those districts where he was on page 1 was 52% while it was 46% in the other districts so that $T = 0.06$

One iteration

Now, drawing a new alphabet (literally calling `sample` in R), we come up with this

UKNRVHPSGQJBODETYWLXIFMCAZ

With this order, Schwarzenegger appears in the 51st position in District 1 -- After rotation, here are the page numbers he appeared on in the 121 districts with multi-page ballots

```
[1] 3 3 3 3 3 3 1 3 3 3 3 3 3 1 5 3 3 3 3 3 2 1 3 2 3 2 3 2 3 1 2 2 1 1 1 1 2 1 2 1  
[38] 2 2 1 1 2 1 2 2 1 2 2 2 2 2 2 1 1 1 1 1 2 2 2 2 2 2 2 1 1 1 1 2 1 1 1 1 2 1 2 1  
[75] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 7 7 7 7 7 7 7 7 2 7 2 7 2 2 7 2 7 2 7 2 7 2 7 7 7  
[112] 6 6 6 6 6 6 6 6 5 6
```

In this case, his average vote share in those districts where he was on page 1 is 44% while it is 50% in the other districts so that $T = -0.06$

One (more) iteration

Drawing yet another new alphabet, we come up with the following

W P B D T X N K L Q E A R J Y G M C O S Z V U I F H

Here, Schwarzenegger appears in the 108th position in District 1 -- After rotation, here are the page numbers he appeared on in the 121 districts with multi-page ballots

```
[1] 5 5 5 5 5 2 5 6 5 5 5 5 2 9 5 5 5 5 5 5 2 6 5 6 5 5 2 5 5 2 2 2 2 4 2 4 1  
[38] 5 4 1 1 4 2 4 5 2 5 5 5 5 4 4 1 1 1 1 1 4 5 4 4 4 5 4 4 2 2 4 2 4 4 1 4 3  
[75] 4 3 4 4 4 3 4 4 4 4 4 4 4 4 4 4 3 3 3 3 3 3 1 3 1 3 1 1 3 1 3 3 1 3 3  
[112] 2 3 2 2 2 2 2 2 2 2 2
```

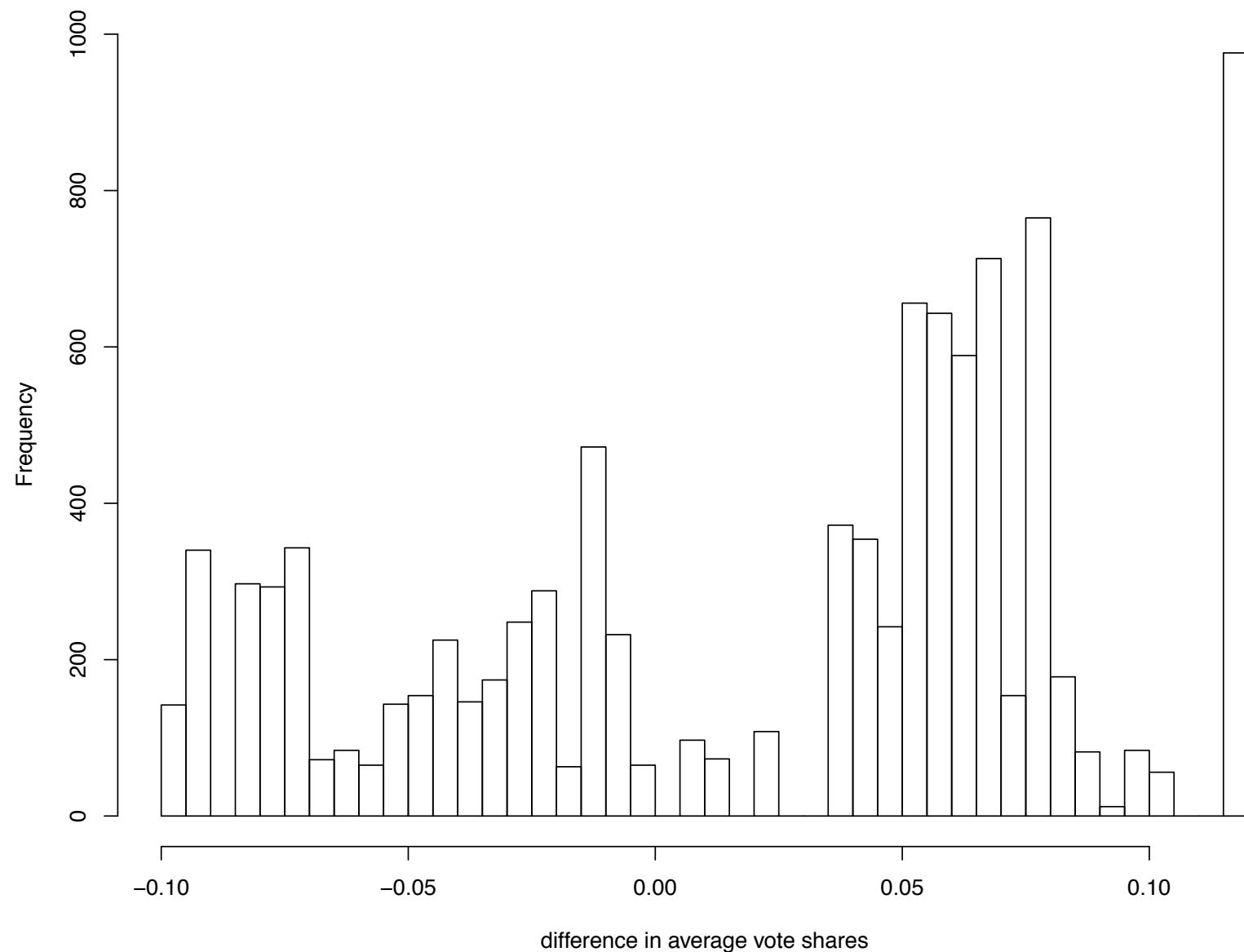
In this case, his average vote share in those districts where he was on page 1 is 55% while it is 47% in the other districts so that **T = 0.08**

Repeating

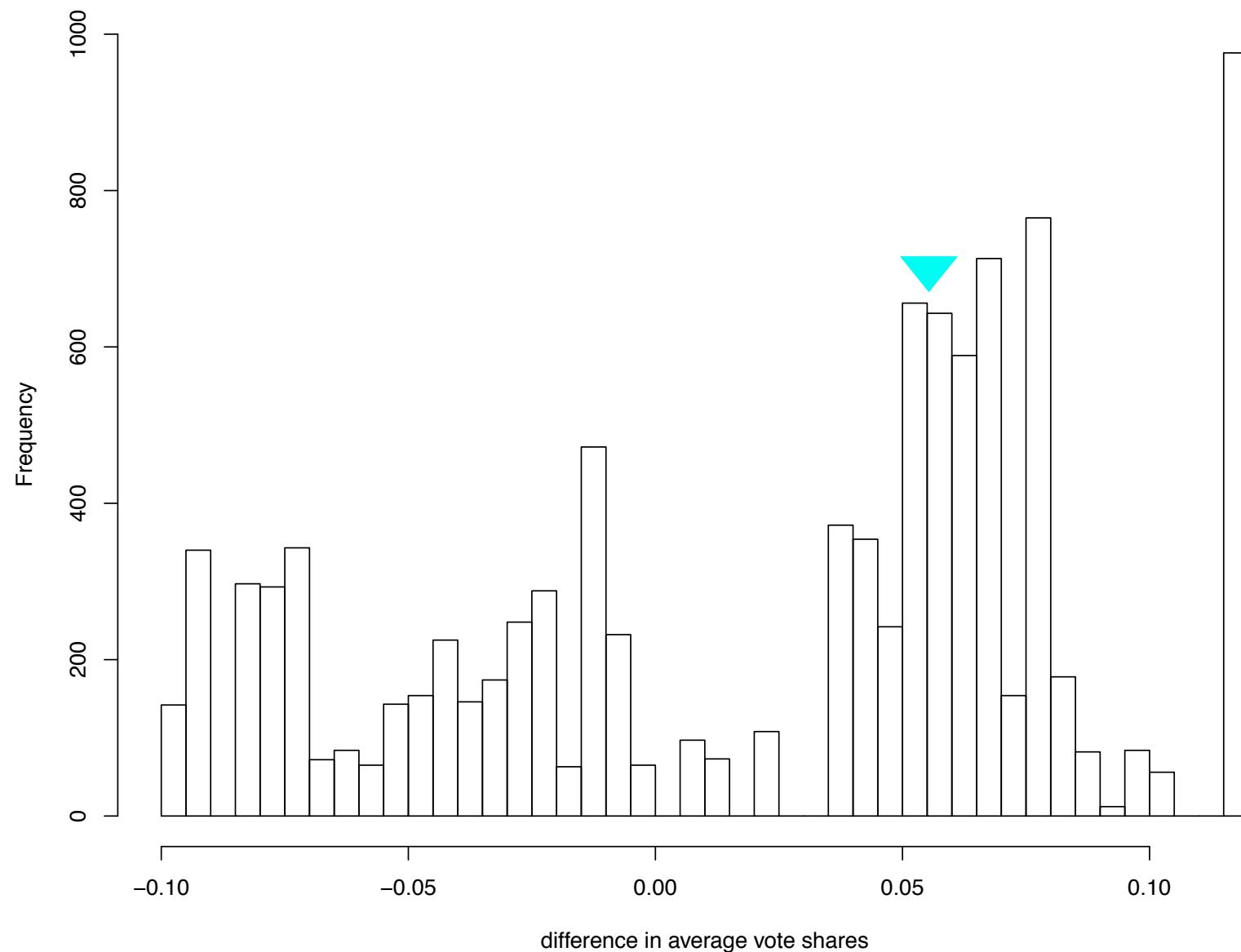
If we carry this out 10,000 times, we come up with **the histogram of T values** on the next page -- These represent the variability in the difference in mean vote shares solely due to the randomization process

Recalling that **the difference for Schwarzenegger in the 2003 election was 0.05** what do you make of this histogram? What can you say about the impact of ballot placement for Schwarzenegger?

Difference in average vote shares, page 1 v. later pages, Schwarzenegger
(10,000 randomized alphabets)



Difference in average vote shares, page 1 v. later pages, Schwarzenegger
(10,000 randomized alphabets)



Formal testing

Ho and Imai computed their P-value as the proportion of re-randomizations that had a difference in average vote shares larger than the observed value of 0.05

In the Schwarzenegger case, 45% of our samples are larger than the observed value of 0.05, meaning he did not see a (statistically) significant boost relative to the uncertainty present due to the randomized alphabets

Formal testing

In technical parlance, this is a “**one-sided test**” meaning that they are looking for evidence against the null hypothesis of 0 difference **only in one tail** (here, the right) of the re-randomization distribution -- This is because prior to conducting their study, they only believed that being on the front page could improve your share of the vote

The authors do comment (but ultimately dismiss) that there might be “regency effects” implying that being listed later in the ballot might actually boost your share -- **If they believed this was possible, then they would look for evidence against the null of 0 difference in both tails**

Formal testing

That is, the P-value would be computed as the proportion of tables with differences that are **either larger than 0.05 and smaller than -0.05** -- You can think of the test statistic now being $|T|$, and to be more extreme here implies large positive or negative differences

Shifting from looking for large values of T to large values of $|T|$ means that (again, before the study was run), the authors had no idea whether the front page listing would increase or decrease a candidate's share and so would take both positive and negative differences from 0 as evidence against the null

Formal testing

In the Neyman-Pearson paradigm, this kind of discussion would boil down to your choice of alternative hypothesis -- With a null that page number has no effect on your vote share, we could entertain three alternatives

A front page listing boosts your share

A front page listing depresses your share

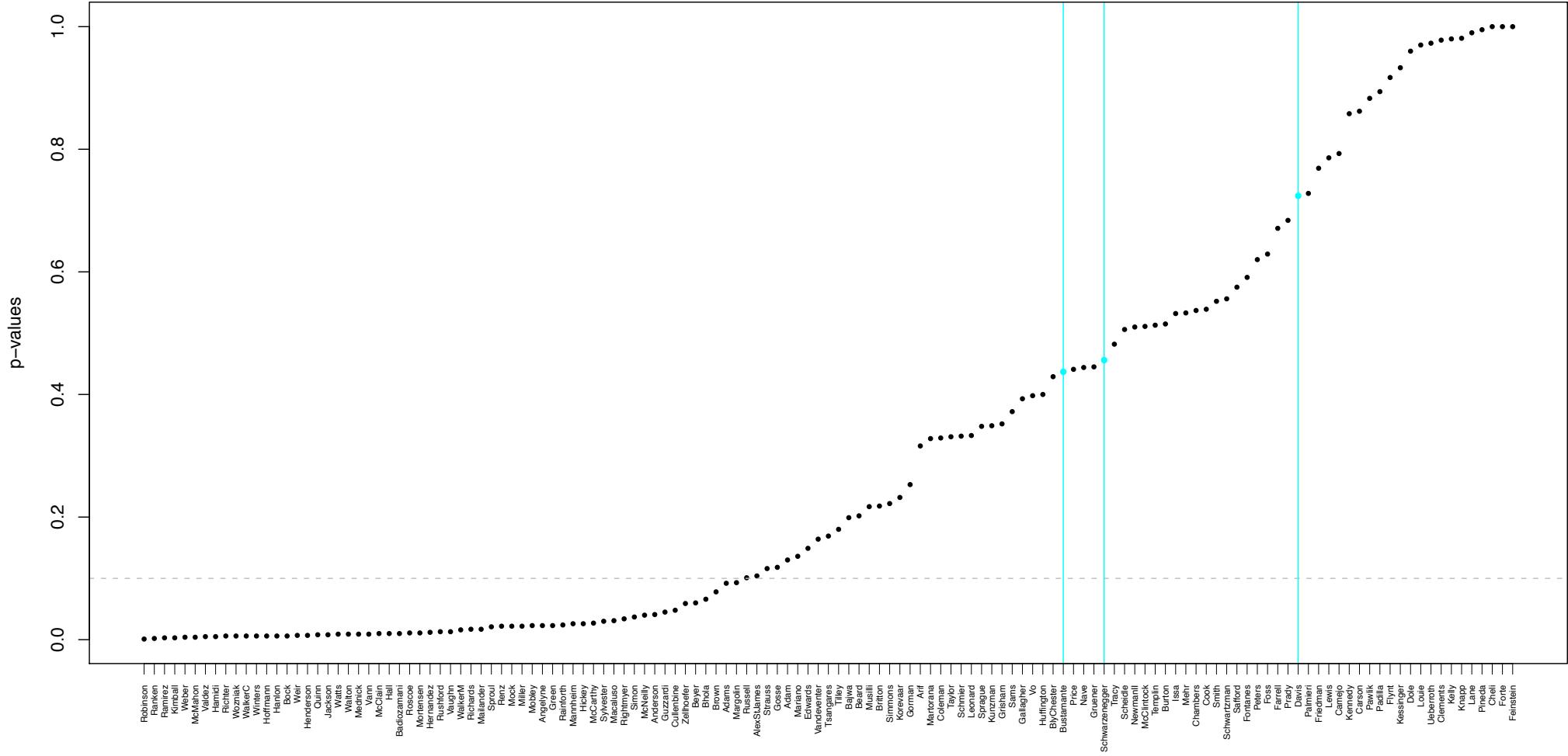
A front page listing affects your share, but we don't know before the study which way the dependence will go

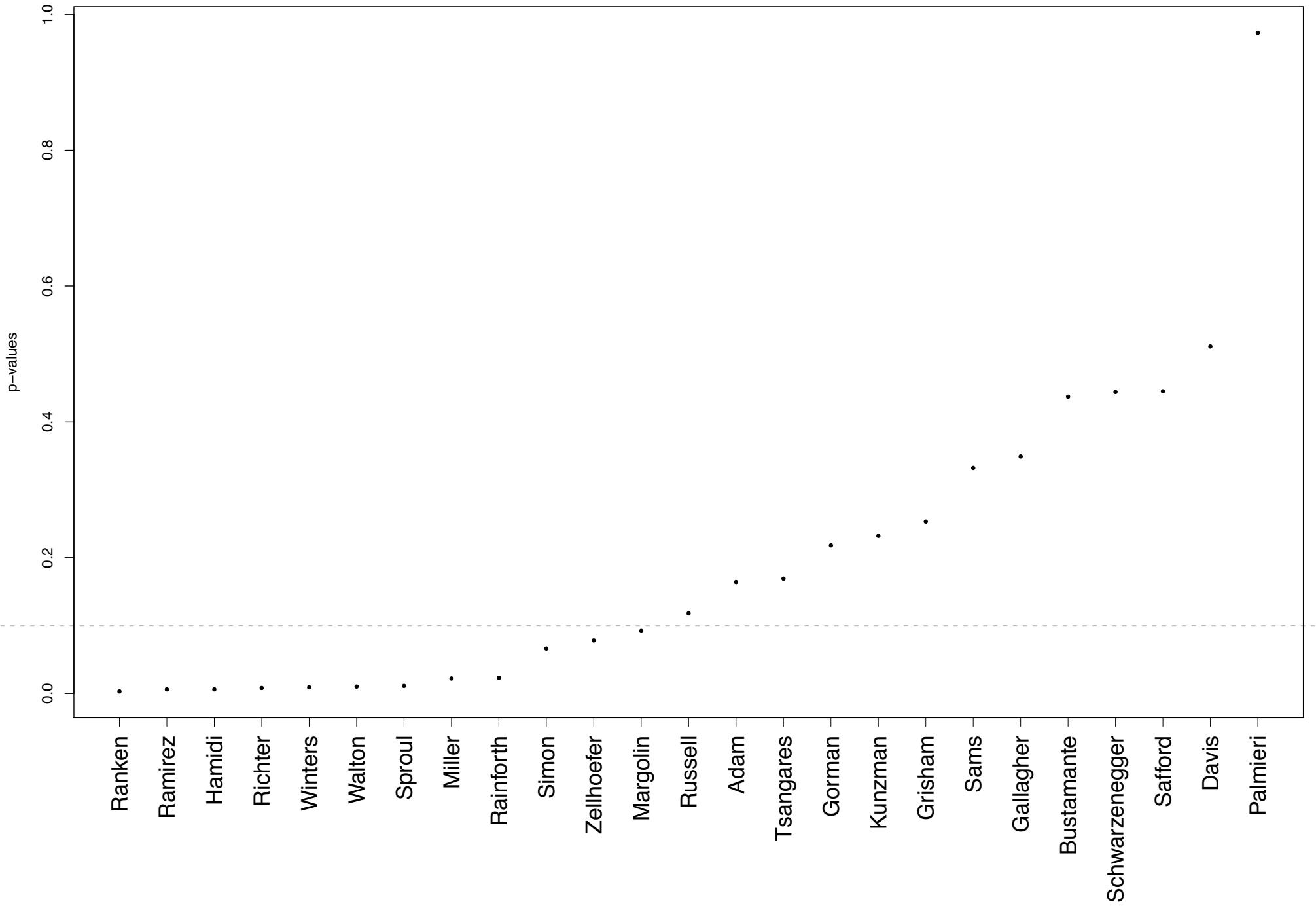
Depending on what you are willing to assume upfront, you will conduct a one-tail test (like the authors do), or a two-tail test (as we suggested with the $|T|$ test statistic)

And the rest of the pack

We can **repeat this analysis** for the remaining 134 candidates from the 2003 ballot and see how they perform -- On the next two pages, we present first the complete list and then a smaller subset (that can actually be read when projected!)

What do you see in these images?





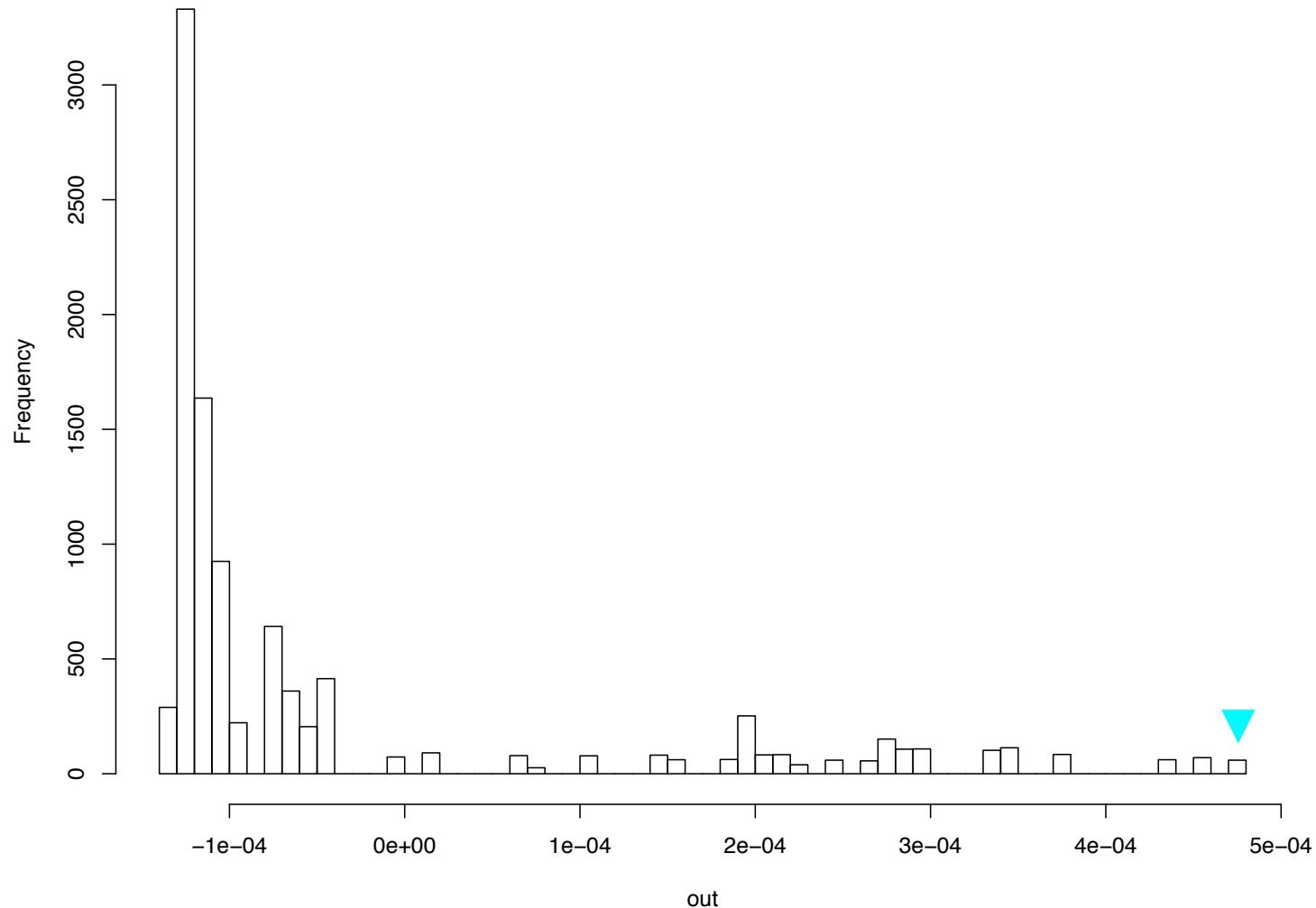
Interpretation

It seems that **for the major candidates, the page-1 effect is undetectable** relative to the uncertainty introduced in the randomization, but that **among minor candidates, we see a fair number of significant effects**

Recall that our testing procedures have a built-in error rate that we can control
-- If we were to reject a null hypothesis of no page effect at the 0.05 level, for example, **we would expect to be mistaken 5% of the time**

in the figure we see 53 out of 135 or **40% of the candidates with P-values smaller than 0.05** -- What this means is that we are seeing more significant effects than we would expect if the null hypothesis was true for all 135 candidates

Difference in average vote shares, page 1 v. later pages, Ranken
(10,000 randomized alphabets)



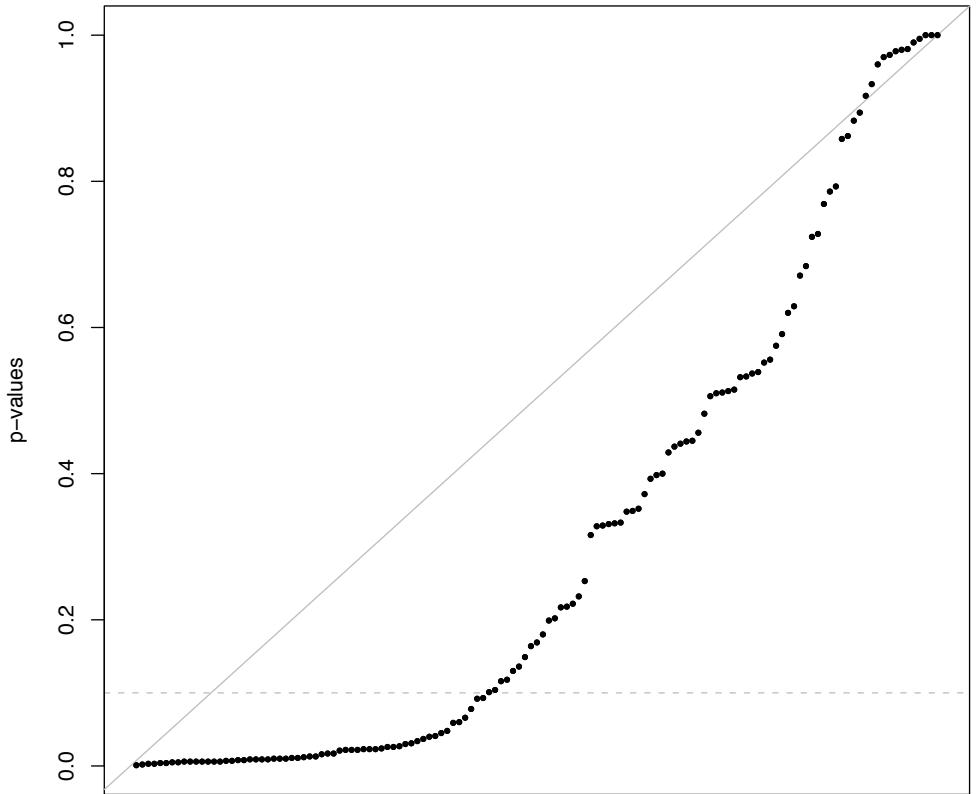
A little advanced

To check that we haven't made a mistake (either in programming or in our thinking), we can try the same experiment but on so-called "**pretreatment variables**" -- These are variables that, unlike vote shares, should **have absolutely no relation to ballot position**

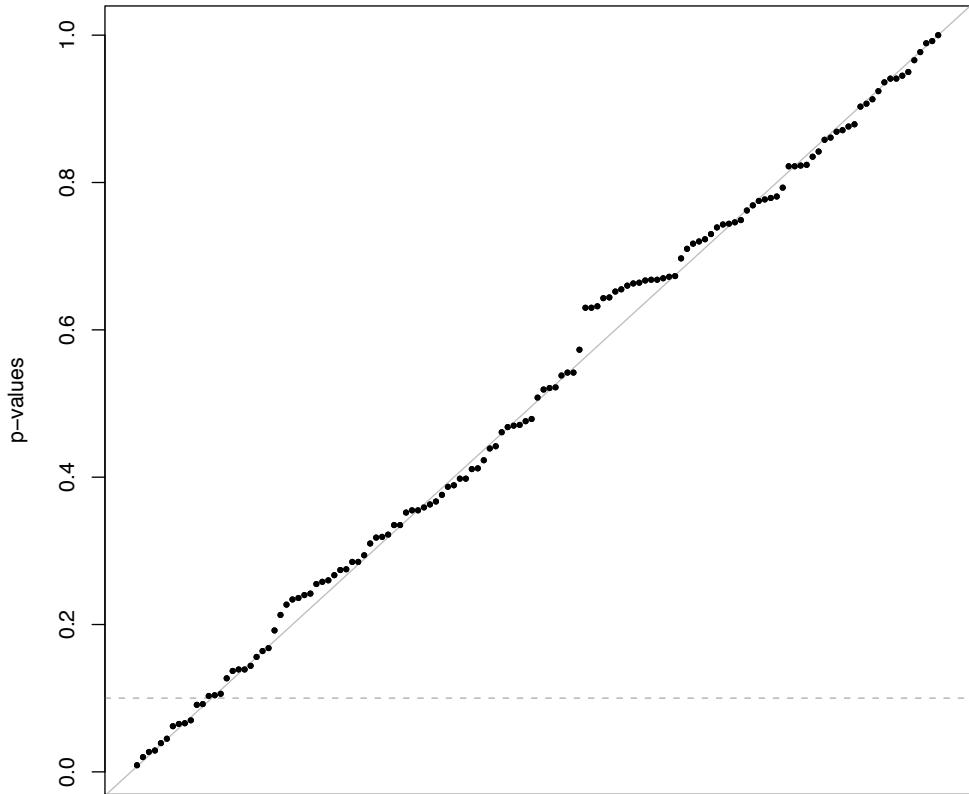
For example, suppose that for each candidate's position relative to our random alphabets, we **compute the difference between average Democratic registration totals for the district-county pairs in which the candidate is listed on page 1 and those for which the candidate is on on a later page**-- This variable had no relationship to the ballot order in the 2003 election and so the null hypothesis is undoubtedly true

Let's look at the sorted P-values for those tests...

p-values testing boost in average vote shared, 135 candidates



p-value plot, green party registration



A little advanced

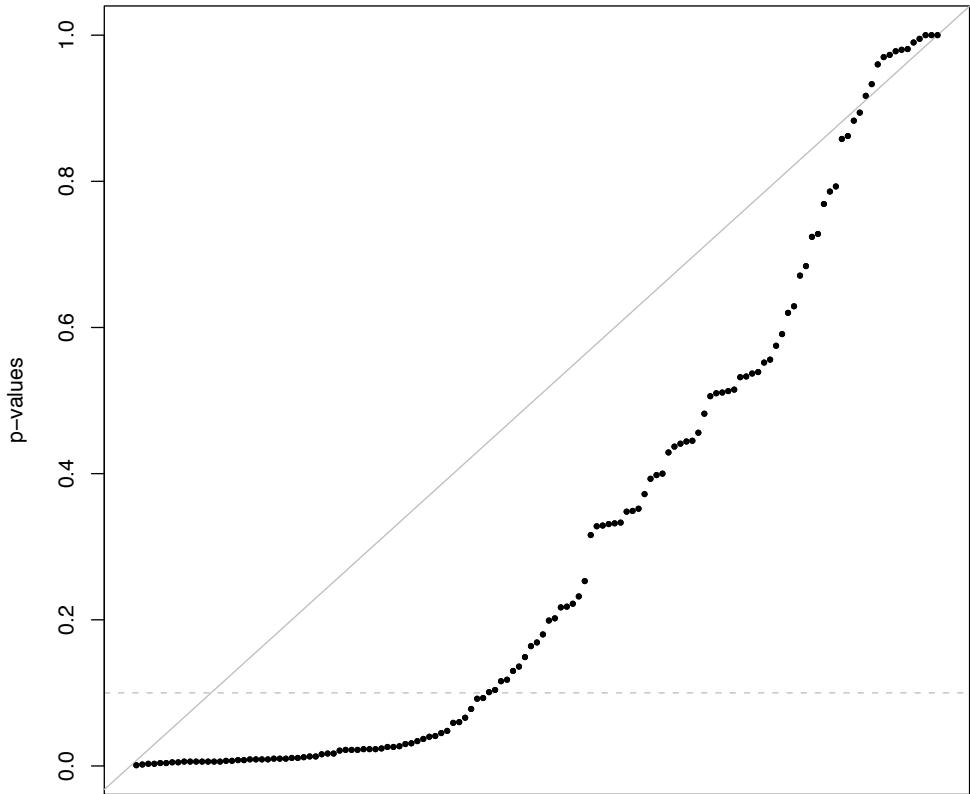
If you think about it a little, if the null hypothesis is true for all 135 candidates, then 5% of P-values should be less than 0.05 -- Similarly, 10% should be less than 0.1 and so on

This explains why, in the previous plot, we looked at **the proportion of candidates** rather than their number and we added **a line with unit slope** -- If all 135 null hypotheses were correct, we'd expect to see our sorted P-values track the line with unit slope

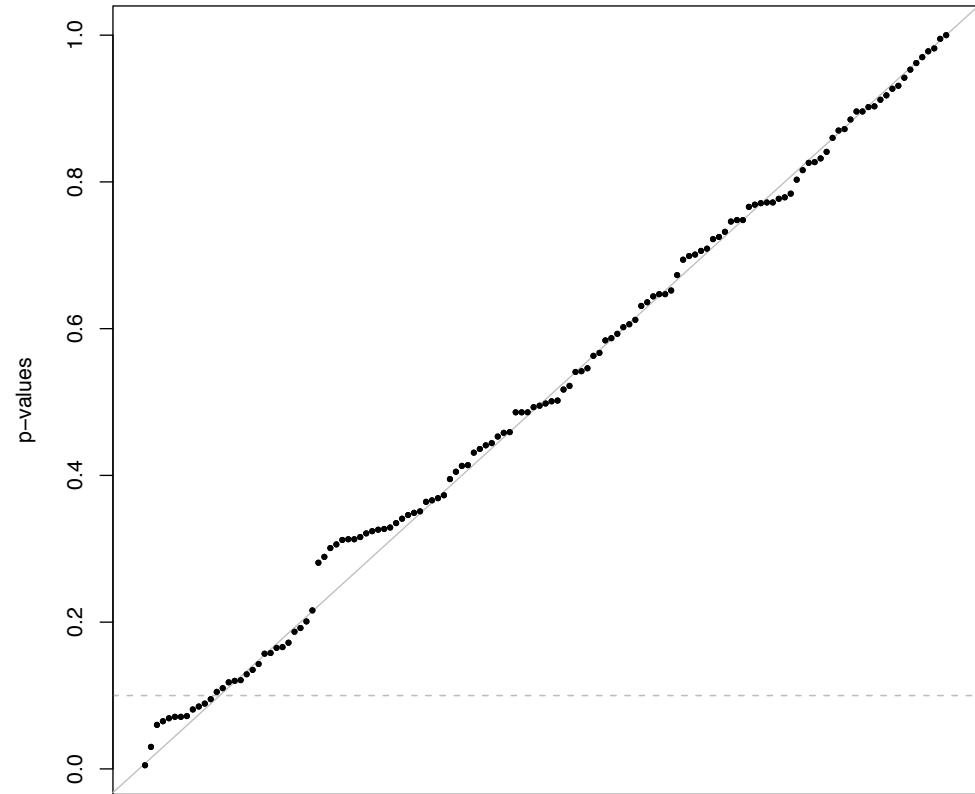
This is the case for Democratic registration totals and (next page) Green Party totals -- In short, **these pictures don't look like the ones we drew for 2003 vote shares**

Based on the P-value plots a few slides back and these “sanity check” plots, Ho and Imai conclude that the **none of the major candidates exhibited significant page-1 effects**, but that about **40% of the minor candidates saw a significant increase in vote shares** by being listed on the first page

p-values testing boost in average vote shared, 135 candidates



p-value plot, democratic registration



To sum

This example illustrates our tag line “**analyze as you randomized**” -- While the California Legislature was not thinking about analysis when they crafted their ballot placement protocol, **it presents us with an interesting “natural experiment”**

The initial randomization opens the possibility of analyzing the results statistically -- We are suddenly able **to make a more refined comment** on the effect of a page-1 spot on a ballot

Ho and Imai take this analysis farther and incorporate “covariates” to adjust for characteristics of the district-county pairs -- This is a bit beyond what we’ve covered so far in class but you can read the paper below if you like

Randomization Inference With Natural Experiments:
An Analysis of Ballot Effects in the 2003 California Recall Election

Daniel E. Ho and Kosuke Imai,
Journal of the American Statistical Association, Vol. 101, No. 475.

Statistical modeling

So far, we have been focusing mainly on re-randomization techniques to assess significance in **study designs that employ some kind of explicit randomization**

In each case, we **defined a test statistic** that represented some aspect of the subjects we were interested in studying and then **created a sampling distribution** for this statistic **under the null hypothesis** that interventions (treatment and control) had no effect on the subjects in our study

The sampling distribution captures the variability present in our experiment under the null hypothesis -- We used this distribution **to judge the size of our observed effect**, deciding whether it was big enough ("extreme enough") to be considered something other than noise

Statistical modeling

The fact that we employed randomization in making our intervention assignments, combined with the null hypothesis of homogeneity between treatment and control provide **a framework for conducting inference**

Random assignments and homogeneity tell us enough about **how the data were generated** (under the null hypothesis) to simulate draws from **the sampling distribution** (under the null hypothesis) and conduct a formal test

These assumptions are relatively weak as statistical assumptions go, and for the next few lectures we will start to add more, fleshing out a framework for **connecting inference (learning from data) to the stochastic (probabilistic) mechanism** that created the data

Modeling

There are many reasons for fitting probability laws to data

We might relate aspects of the distribution (features or parameters) to a scientific theory that reveals something about the state of Nature

Probability models can be used for purely descriptive purposes, acting as a kind of data summary or “compression”

Finally, we are often interested in simulation, in using these models to make predictions or to generate new data that can be fed into a larger modeling exercise

The normal family

The normal or Gaussian distribution involves two parameters, μ and σ , where μ is the mean (the center) and σ is the standard deviation (the spread)

$$\mathcal{F} = \left\{ f(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/2\sigma^2} : \mu \in \mathbb{R}, \sigma^2 > 0 \right\}$$

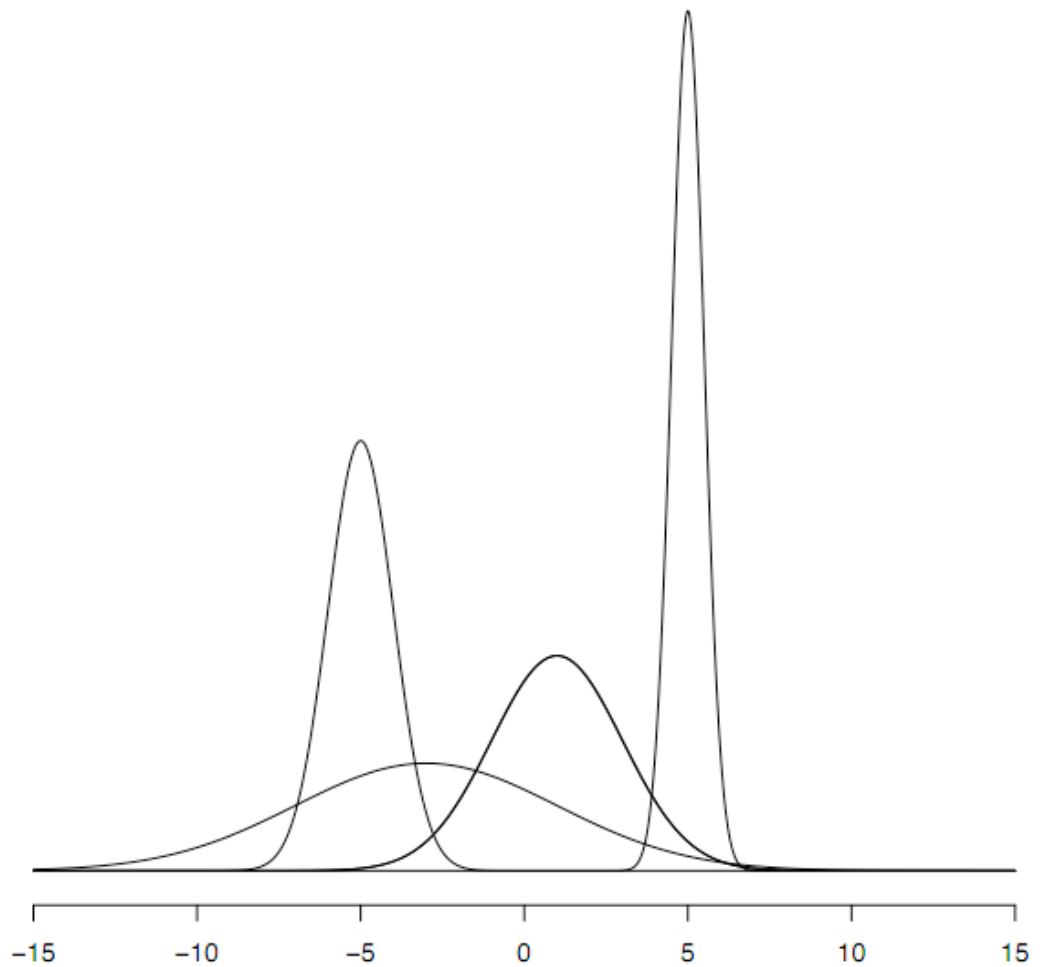
The normal family

The normal family is often used as an “error distribution,” representing a large number of independent random fluctuations*

We recall that we can generate any member of this family from the so-called standard normal distribution -- That is, if Z has a normal distribution with mean zero and standard deviation one, then

$$X = \mu + \sigma Z$$

again has a normal distribution but with mean μ and standard deviation σ

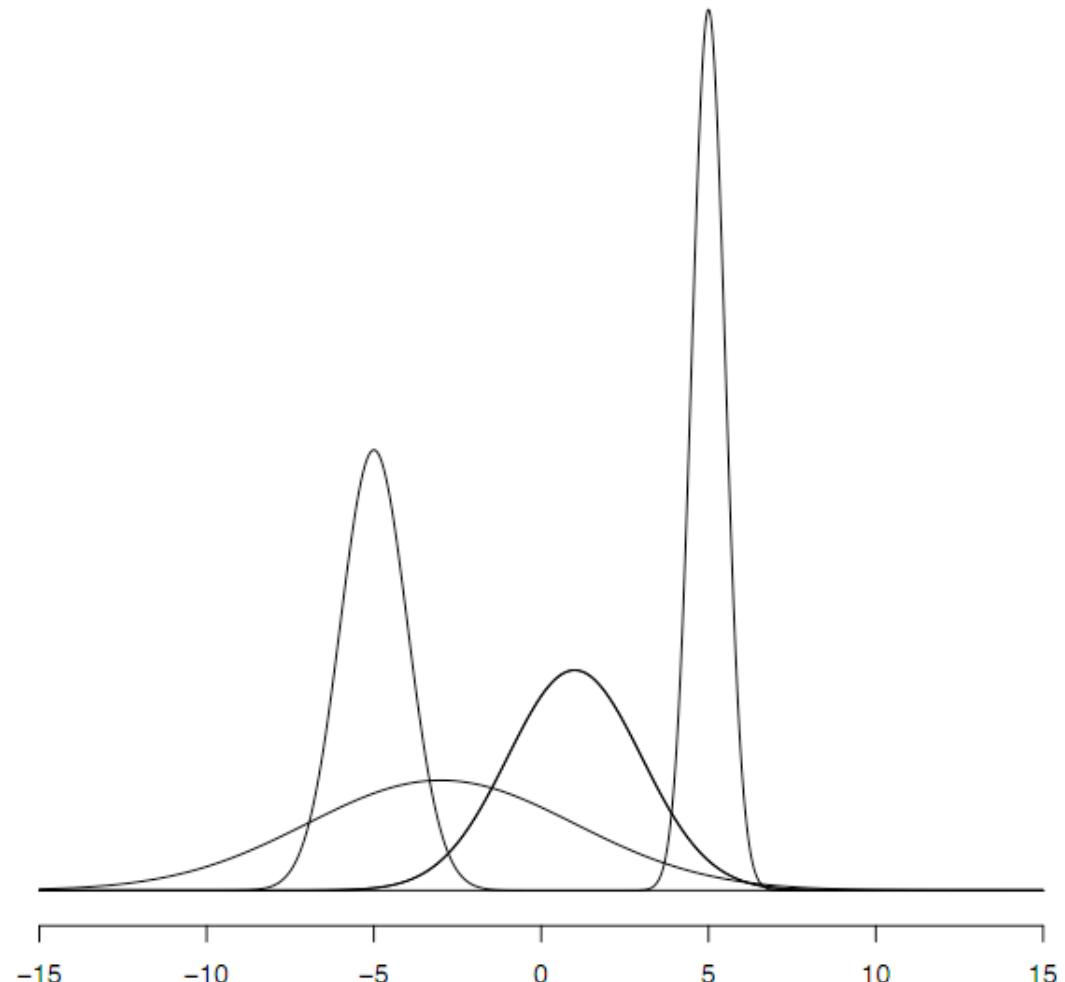


The normal family

Put another way, this is a **location-scale family** meaning that every member of it can be written as

$$f(x|\mu, \sigma) = \frac{1}{\sigma} f_0\left(\frac{x - \mu}{\sigma}\right)$$

where $f_0(x) = e^{-x^2/2}/\sqrt{2\pi}$ is the standard normal density



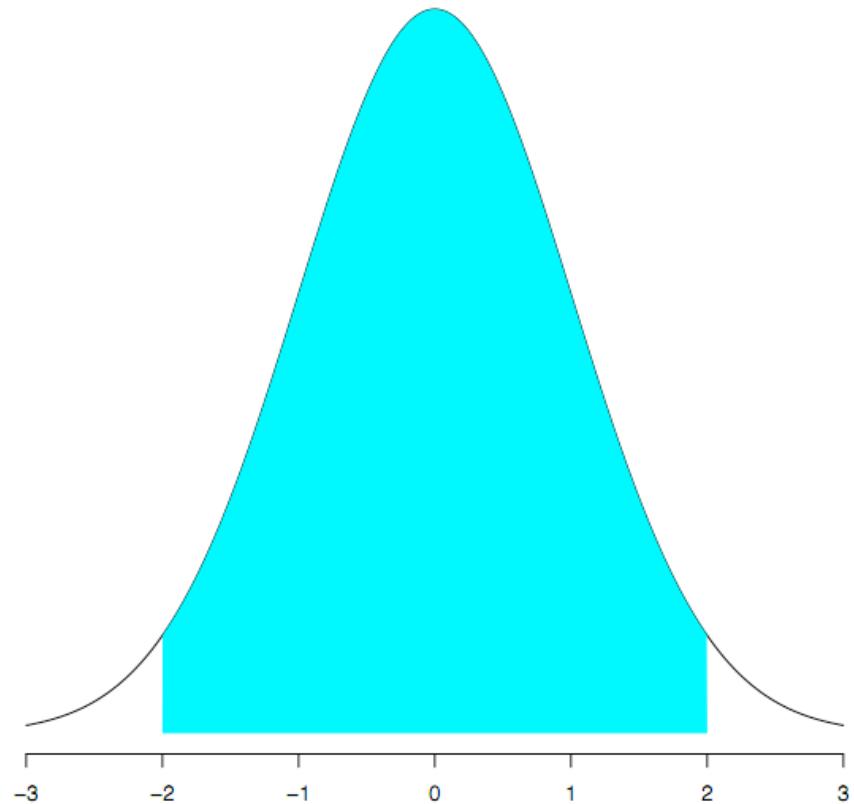
The normal family

The location-scale property means that we can restate questions about **arbitrary normal random variables** in terms of a **standard normal**

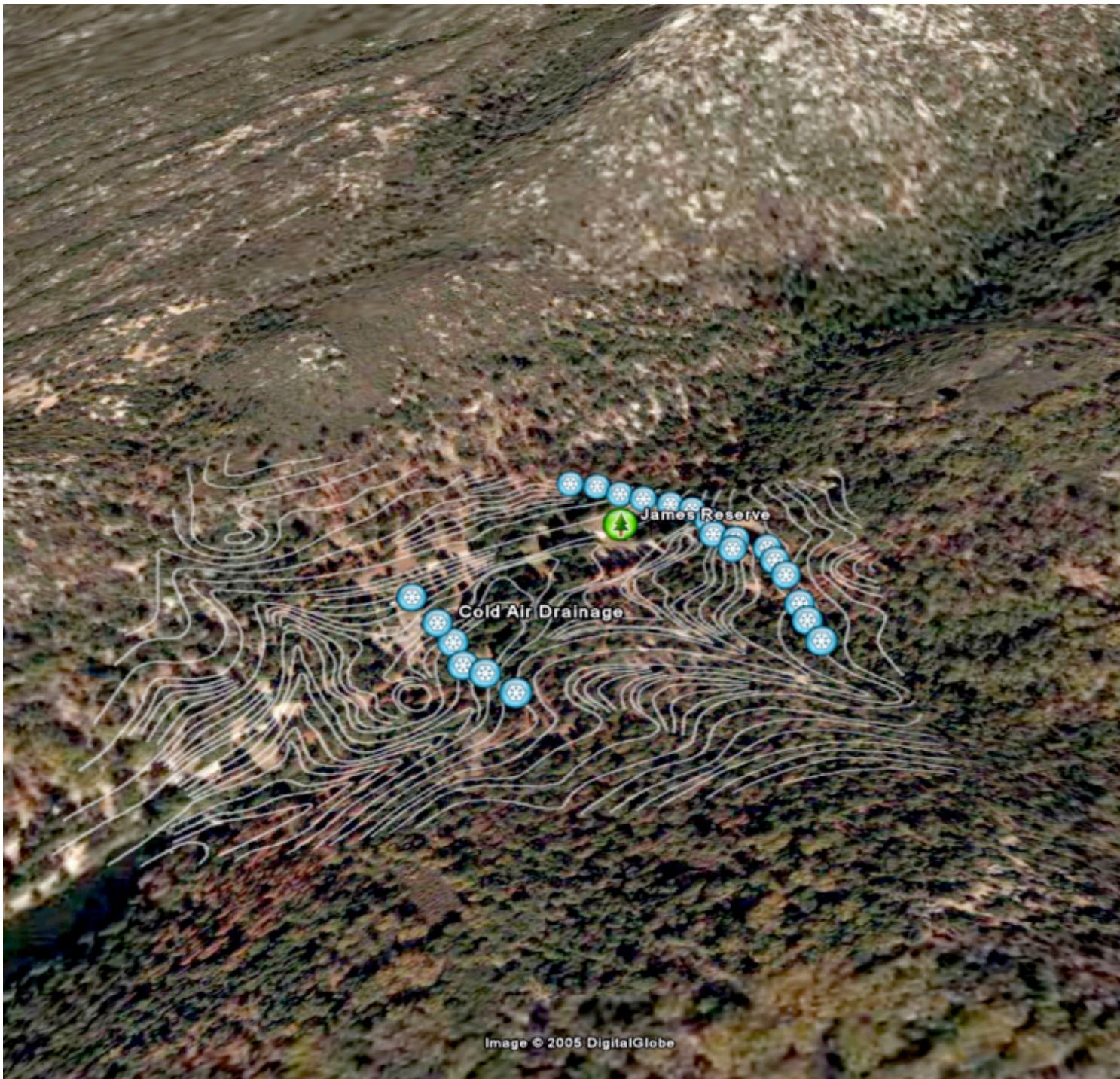
For example, from the fact that

$$P(-1.96 \leq Z \leq 1.96) = 0.95$$

where Z is a standard normal, we can say that any normal random variable has about a 95% chance of being within two standard deviations of its mean





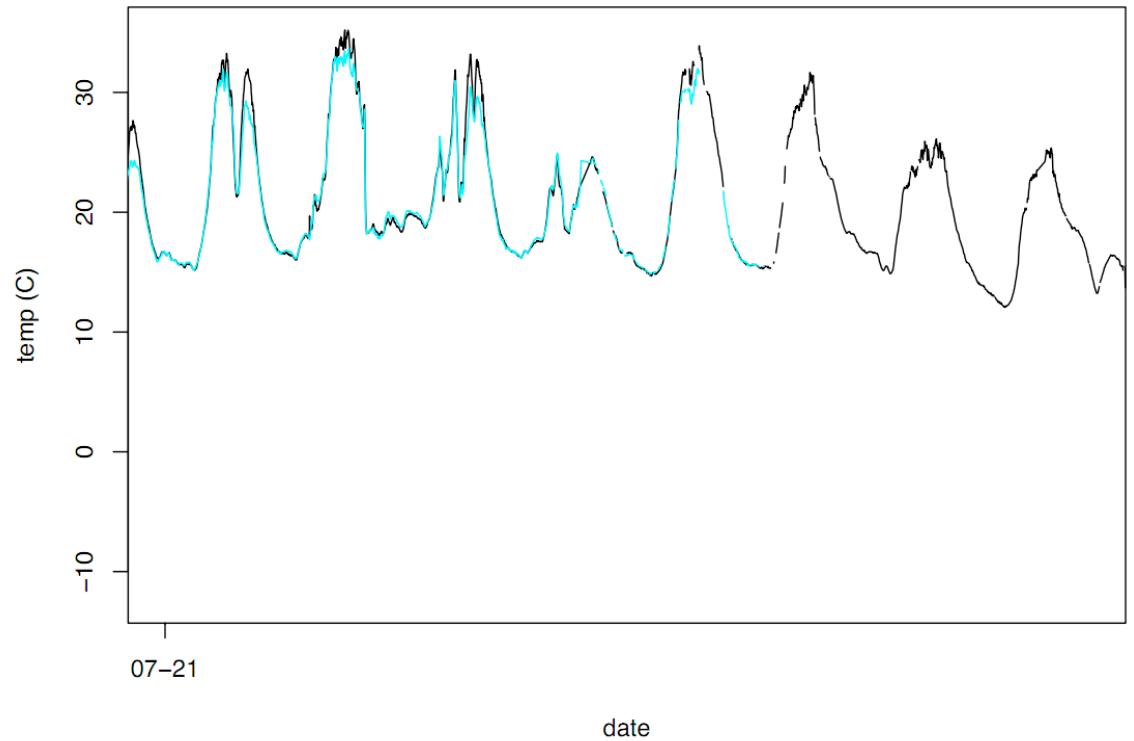


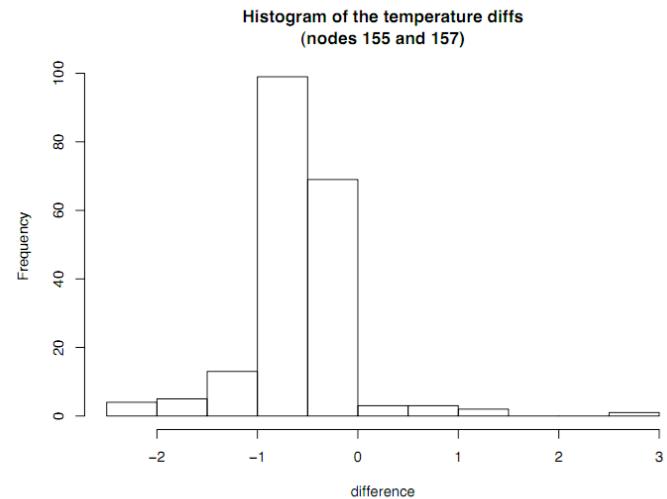
Example: CAD

Here are two sets of temperature readings from neighboring sensors in the cold air drainage transect at the James Reserve

Let's consider the difference between two neighboring sensors in the transect -- In principle, these measurements should be similar

We will initially consider a simple normal model for the differences at a given time of day, say 11a

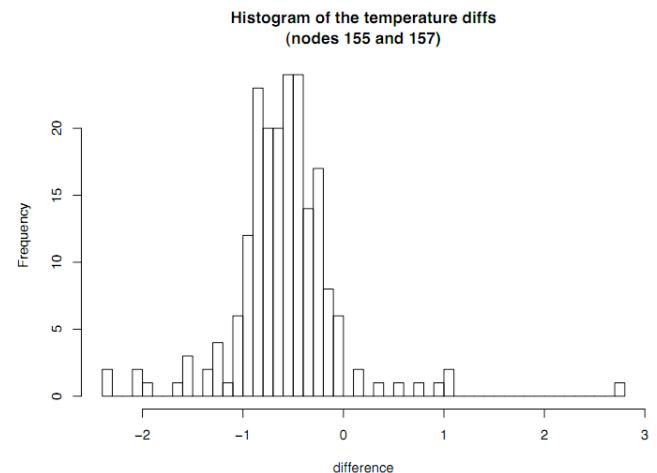
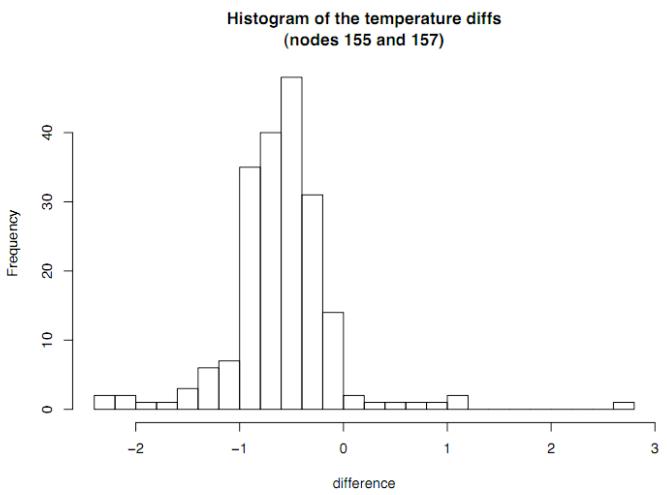




Example: CAD

At the right we have a histogram of the difference in temperatures for nodes 156 and 157 at 11am for 225 days

Here we are varying the bin counts to give us a better sense of the shape of the distribution -- Any comments?



The binomial family

Let X denote the number of successes in m independent trials, each with success probability p -- Then the probability function of X belongs to the family

$$\mathcal{F} = \left\{ f(k|p) = \binom{m}{k} p^k (1-p)^{m-k} : p \in [0, 1] \right\}$$

Unlike the normal, the binomial is indexed by **a single parameter, p** , the success probability (we consider m , the number of trials, to be fixed and given)

Recall that if X has a binomial distribution (m,p) , then **X has expected value mp and variance $mp(1-p)$**



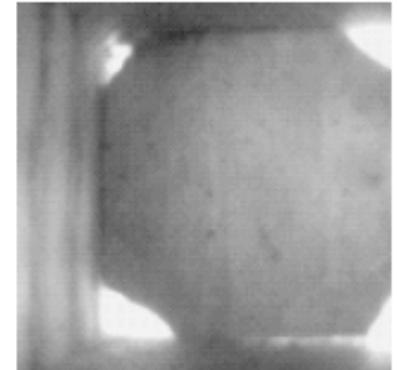


Nestbox 8

Images were taken of the inside of the nest box every 90 minutes for a period of two weeks (14 days) -- The images were manually tagged to indicate the presence or absence of an (adult) bird

For each 90 minute period, we can think of the $m=14$ trials as tossing a coin with probability p that we see a bird -- We then might think of the 15 counts as being observations of binomial random variables

	d1	d2	d3	d4	d5	d6	d7	d8	d9	d10	d11	d12	d13	d14
1	0	1	1	1	1	1	1	1	1	0	1	1	0	1
2	1	1	0	1	0	1	0	1	1	1	0	1	1	1
3	0	0	0	1	0	0	0	1	1	1	0	1	1	1
4	0	1	1	0	1	0	0	0	1	0	1	0	1	1
5	1	1	0	1	0	1	0	1	1	0	0	1	1	1
6	1	1	0	1	1	1	0	0	1	1	0	1	1	1
7	1	1	0	1	1	1	0	1	1	0	0	1	1	1
8	1	1	0	1	1	1	0	1	1	1	0	1	1	1
9	1	1	1	0	1	1	1	1	1	1	0	1	1	1
10	1	1	0	0	1	1	0	1	1	0	1	1	1	1
11	0	1	1	0	0	1	1	1	0	0	1	0	1	1
12	0	1	0	0	1	1	1	0	1	0	1	0	1	1
13	0	0	0	0	0	0	0	0	0	0	0	0	0	0
14	1	1	0	1	0	1	0	1	1	0	0	1	0	1
15	1	1	0	1	1	1	1	1	1	1	0	1	1	1



Point estimation

A point estimate refers to the assignment of a single value to some quantity of interest based on data -- **In a parametric model**, for example, we want to identify or **learn sensible values for the parameters**

Estimates earn themselves hats -- So an estimate of μ is denoted $\hat{\mu}$, $\hat{\sigma}$ for σ and so on

Point estimation

To set notation a little, we will assume that X_1, \dots, X_n are a sample of n independent observations drawn from $f(x|\theta^*)$, a member of the parametric family of probability functions $f(x|\theta)$ indexed by a (possibly vector valued) parameter θ

For the next few slides, we will be interested in forming an estimate of θ^* based on data -- That is, we will examine ways to form $\hat{\theta} = \hat{\theta}(X_1, \dots, X_n)$ using the observations X_1, \dots, X_n

Maximum likelihood

Assuming x_1, \dots, x_n are **independent and all drawn from the same parametric probability function** $f(x|\theta)$ then their joint distribution is given by

$$f(x_1, \dots, x_n | \theta) = \prod_{i=1}^n f(x_i | \theta)$$

We now view this as a function of θ rather than the x_i -- Specifically, we define **the likelihood function** to be

$$\mathcal{L}(\theta) = \prod_{i=1}^n f(X_i | \theta)$$

with **the log-likelihood function**

$$l(\theta) = \log \mathcal{L}(\theta) = \sum_{i=1}^n \log f(x_i | \theta)$$

Maximum likelihood

As its name suggests, as an estimation procedure, maximum likelihood suggests selecting a value for $\hat{\theta}$ that makes the data the most likely or probable

Formally, this means $\hat{\theta} = \operatorname{argmax} \mathcal{L}(\theta)$

Maximum likelihood

The first example of this approach is due to Daniel Bernoulli who (in 1777) published a paper on the treatment of errors in which he proposed the following

Of all the innumerable ways of dealing with errors of observation, one should choose the one which has the highest degree of probability for the complex of observations as a whole.



Maximum likelihood

R. A. Fisher is really responsible for establishing many of the properties we now associate with the method

What we can find from a sample is the likelihood of any particular value of [a parameter] ρ , if we define the likelihood as the quantity proportional to the probability that, from a population having that particular value of ρ , a sample having the observed value ρ should be obtained. So defined, the probability and likelihood are quantities of an entirely different nature.



The normal family

Let's derive the maximum likelihood estimates for observations coming from the normal family -- We can rewrite the likelihood slightly

$$\begin{aligned}\mathcal{L}(\mu, \sigma) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{(X_i - \mu)^2 / 2\sigma^2} \\ &\propto \sigma^{-n} \exp \left[-\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2 \right] \\ &= \sigma^{-n} \exp \left[-\frac{nS^2}{2\sigma^2} \right] \exp \left[-\frac{n(\bar{X} - \mu)^2}{2\sigma^2} \right]\end{aligned}$$

where we have inserted the sample mean and (nearly) the sample variance

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \quad \text{and} \quad S^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

The normal family

In terms of the log-likelihood function, this becomes (up to a constant C that doesn't depend on μ or σ)

$$l(\mu, \sigma) = -n \log \sigma - \frac{nS^2}{2\sigma^2} - \frac{n(\bar{X} - \mu)^2}{2\sigma^2} + C$$

Now, we take partial derivatives with respect to the parameters μ and σ to find the maximum likelihood estimates (MLEs) $\hat{\mu} = \bar{X}$ and $\hat{\sigma} = S$

The normal family

I have to admit that running through that (albeit not that horrible math) and coming up with the sample mean and standard deviation as the MLEs is **a little bit of a let-down**

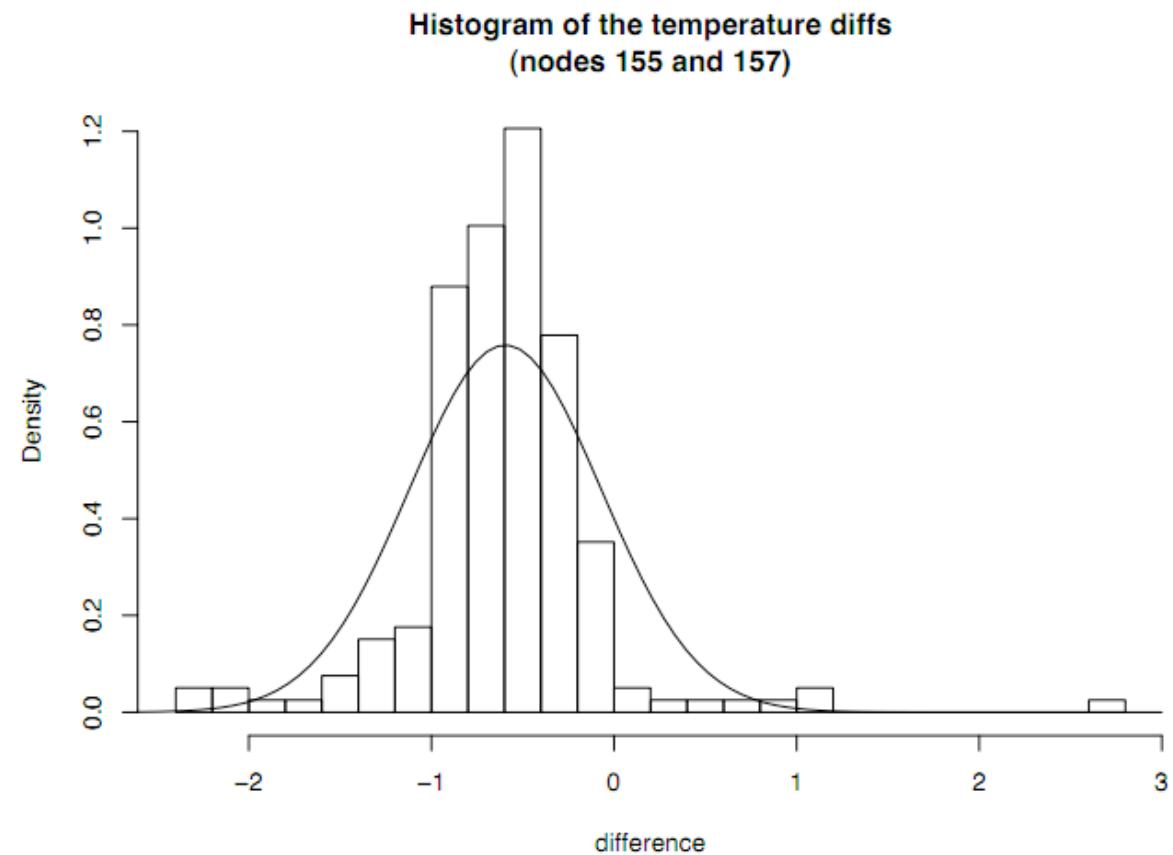
Granted, the estimates make intuitive sense and as we will see, being an MLE provides you with certain properties you might not otherwise expect

However, this does beg the question, if all we're doing here is computing the sample mean and variance, well, we could do that for any set of data -- **Fitting a model doesn't mean it, um, fits**

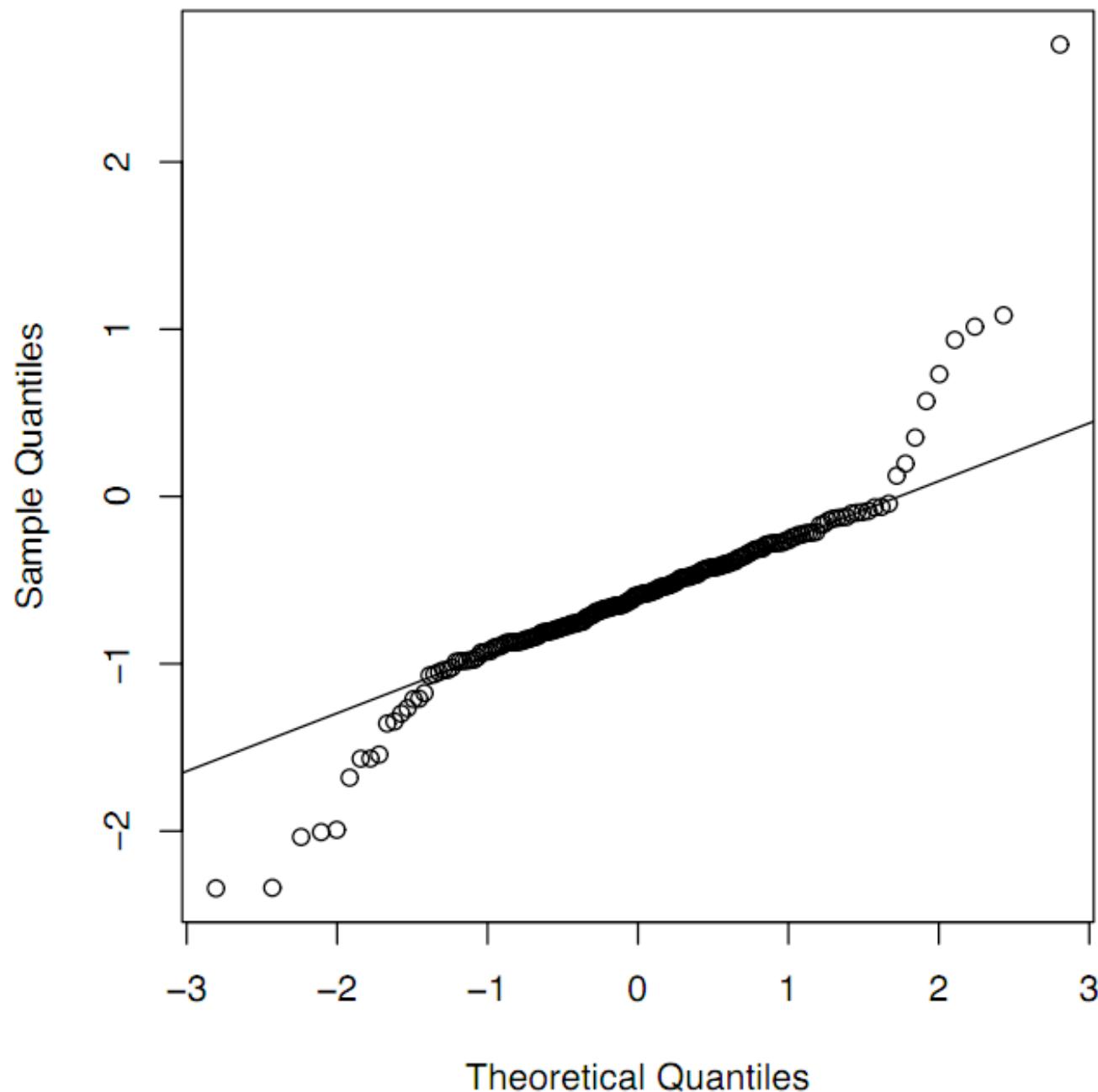
Example: CAD

At the right we have one of our histograms of the difference data with an overlay of our fitted normal (here, $\hat{\mu} = -0.6$ and $\hat{\sigma} = 0.53$)

How well have we done? Is this the best display to assess the fit?



Normal Q-Q Plot



The binomial family

Let's apply ML to the binomial family -- Specifically we have a series of binomial observations which give rise to the likelihood

$$\begin{aligned}\mathcal{L}(p) &= \prod_{i=1}^n p^{X_i} (1-p)^{m-X_i} \\ &= p^U (1-p)^{nm-U}\end{aligned}$$

where we define $U = X_1 + \dots + X_n$

This gives us a log-likelihood of the form

$$l(p) = U \log p + (nm - U) \log (1-p) + C$$

and after differentiating with respect to p , we find the MLE to be $\hat{p} = U/nm$

Nestbox 8

Doing a simple row sum on our data table, we have

11 10 7 7 9 10 10 11 12 10 8 8 0 8 12

that we postulate as independent observations from a binomial distribution with
 $m=15$ -- Our MLE is $\hat{p} = 133/210 = 0.633$

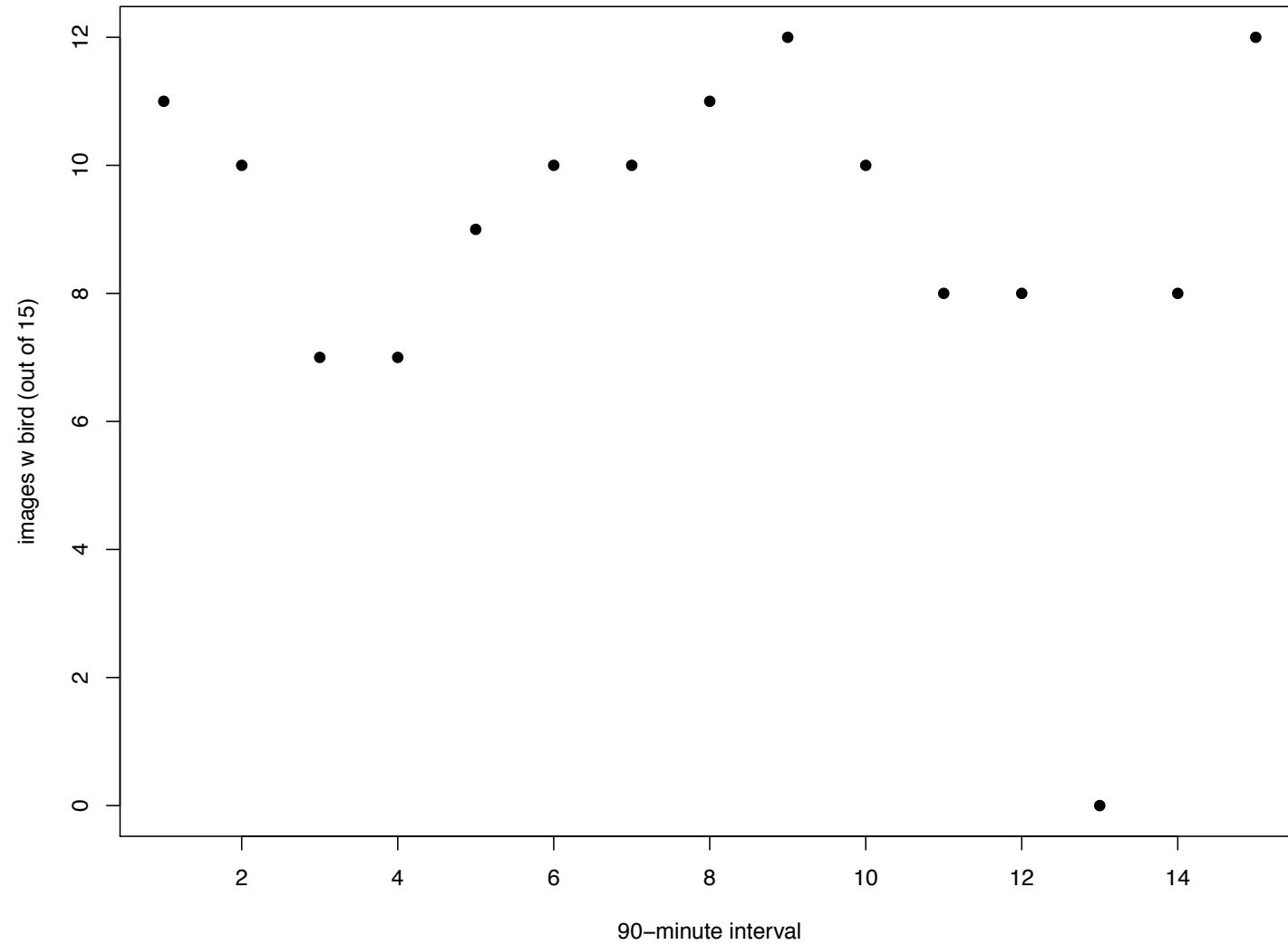
What kind of model check might we perform here?

Nestbox 8

Notice that the binomial is different from the normal in that it only has one parameter that feeds both the mean and variance

In these cases, we often check for what is known as overdispersion -- That is, is the sample variance is a lot bigger than what the model would predict

For example, the sample variance of our 15 numbers is 8.7 but the model predicts $15 \cdot 0.63 \cdot 0.37 = 3.5$

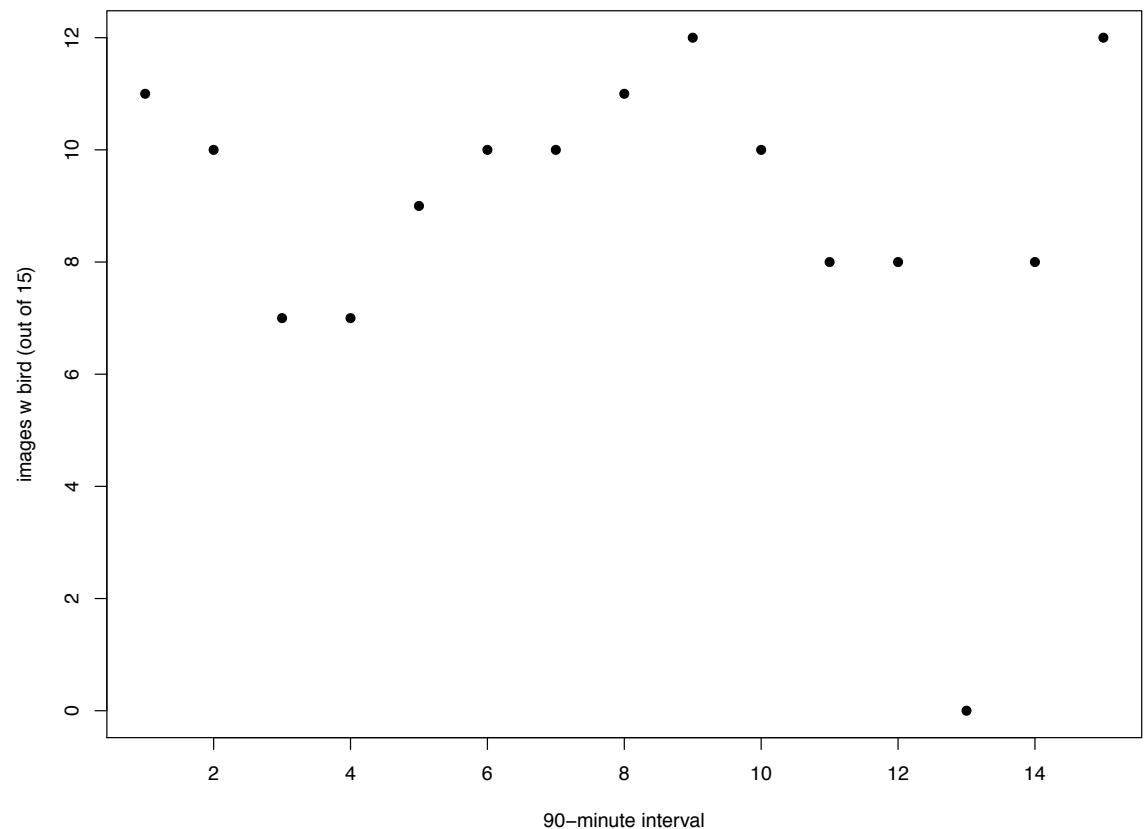


Nestbox 8

If we look at the data we see that there was one 90-minute interval when the bird consistently failed to make an appearance

That zero is not well captured by a binomial (the chance of getting a zero under the fitted model is $0.37^{15} = 0.0000003$) -- When we remove the 0, the variance drops to a more reasonable 2.88

There are various ways to improve a binomial model -- We might, for example think about a way to add extra zeroes



Nestbox 8

The zero, together with the apparent smoothness (across 90-minute intervals) of our counts suggests that it might be unrealistic to assume the p is the same for each time period

Instead, we might posit a model in which p varies smoothly over time, with larger values at night and midday

