

Bayes rule

- recall def of conditional:

- ▶ $P(a|b) = P(a \wedge b) / P(b)$ if $P(b) \neq 0$

$$P(a|b)P(b) = P(a \wedge b) = P(b|a)P(a)$$

$$P(a|b) = P(b|a)P(a) / P(b) \quad P(b) \neq 0$$

Geoff Gordon—10-701 Machine Learning—Fall 2013

mult thru by $P(b)$: $P(a|b)P(b) = P(a \wedge b)$

holds even if $P(b)=0$ (check), so some take this as basic def of conditioning

so: $P(a \wedge b) = P(a|b)P(b) = P(b|a)P(a)$

Bayes rule: divide last eq by $P(b)$ (if $P(b) \neq 0$)

$$P(a|b) = P(b|a)P(a) / P(b)$$

note: if b is observed, don't need to worry about $P(b) = 0$

good rule: only condition on events we might observe :-)

Bayes rule w/ background event C (just like changing to a smaller universe)

$$P(a \wedge b|C) = P(a|b,C)P(b|C) = P(b|a,C)P(a|C)$$

Bayes rule: sum version

- $P(a | b) = P(b | a) P(a) / P(b)$

Geoff Gordon—10-701 Machine Learning—Fall 2013

what if we don't know $P(b)$, but still know $P(b|a)$ and $P(a)$
(often happens)
suppose MEEP $A = a_1, a_2, \dots, a_n$ for moderate n

we know $\sum_i P(a_i | b) = 1$
so $\sum_i P(a_i | b) P(b) = P(b)$ (mult by $P(b)$ on both sides)
LHS is $\sum_i P(b|a_i) P(a_i)$ (another application of Bayes rule)
each term assumed known, sum tractable since n moderate

Bayes rule in ML

- $P(\text{model} \mid \text{data}) = P(\text{data} \mid \text{model}) P(\text{model}) / P(\text{data})$

Geoff Gordon—10-701 Machine Learning—Fall 2013

3

Why do we care about Bayes? Lets us take info about conditional in one direction ($P(\text{data} \mid \text{model})$) and get info about the other direction ($P(\text{model} \mid \text{data})$).

LHS $P(\text{model} \mid \text{data})$ tells us which models are probable give the data. This is (mostly) what we want out of ML! Called “posterior”

$P(\text{data} \mid \text{model})$ usually has an explicit formula, e.g., gaussian distribution: $P(x_i \mid \mu) = \exp(-(x_i - \mu)^2 / 2) / \sqrt{2\pi}$. Called “likelihood”

$P(\text{model})$ = “prior” -- sometimes controversial but not hard to write a minimally reasonable one

$P(\text{data})$: “normalizing constant” (also “annoyance”) -- often hard to get

can use sum version of Bayes rule (sum numerator over models)

OK if not too many models under consideration

can use sum version + approximation tricks (numerical quadrature, MCMC like Gibbs)

another idea on next slide

note: normalizing constants are often ignored. This is a common pattern, but it doesn't mean it's always safe (or a good idea) to ignore them! Sometimes all of the useful information is in the normalizing constant...

Bayes rule vs. MAP vs. MLE

- $P(\text{model} \mid \text{data}) = P(\text{data} \mid \text{model}) P(\text{model}) / P(\text{data})$

Geoff Gordon—10-701 Machine Learning—Fall 2013

4

MAP: just take model for which RHS is highest -- if we have to choose just one point from posterior density, why not this one?

[sketch]

now $P(\text{data})$ really is ignorable (doesn't change order)

Seems like a horrible approximation (from Bayesian perspective), but actually has theory to back it up.

MLE: ignore $P(\text{model})$ term too. Why? Because people argue about the prior; because with enough evidence from data it might become negligible.

Seems even worse, but still has theory to back it up. (see next slide)

Frequentist vs. Bayes



Jerzy Neyman

FIGHT!!!



rev. Thomas Bayes

- Nature as adversary vs. Nature as probability distribution
- Probability as long-run frequency of repeatable events vs. odds for bets I'm willing to take

Geoff Gordon—10-701 Machine Learning—Fall 2013

treat Nature as a probability distribution vs. an adversary
Bayes: if you have the right distribution over models and can compute with it, good things happen
Frequentist: but both of those are horrible assumptions
Plenty of paradoxes for both sides if you want to cast stones

Test for a rare disease

- About 0.1% of all people are infected
- Test detects **all** infections
- Test is highly specific: 1% false positive
- You test positive. What is the probability you have the disease?

Geoff Gordon—10-701 Machine Learning—Fall 2013

P(disease)
P(test | +disease)
P(test | -disease)

P(+disease | +test)
= P(+test | +disease) P(+disease) / P(+test)
P(+t) = P(+t|+d)P(+d) + P(+t|-d)P(-d)
= 1 * .001 / (1*.001 + .01*.999)
= .091

Test for a rare disease

- About 0.1% of all people are infected

Bonus: what is probability an average med student gets this question wrong?

the disease?

Geoff Gordon—10-701 Machine Learning—Fall 2013

$P(\text{disease})$

$P(\text{test} \mid +\text{disease})$

$P(\text{test} \mid -\text{disease})$

$P(+\text{disease} \mid +\text{test})$

$= P(+\text{test} \mid +\text{disease}) P(+\text{disease}) / P(+\text{test})$

$P(+\text{test}) = P(+\text{test} \mid +\text{disease}) P(+\text{disease}) + P(+\text{test} \mid -\text{disease}) P(-\text{disease})$

$= 1 * .001 / (1 * .001 + .01 * .999)$

$= .091$

Follow-up test

- Test 2: detects 90% of infections, 5% false positives
 - ▶ $P(+\text{disease} \mid +\text{test1}, +\text{test2}) =$

Geoff Gordon—10-701 Machine Learning—Fall 2013

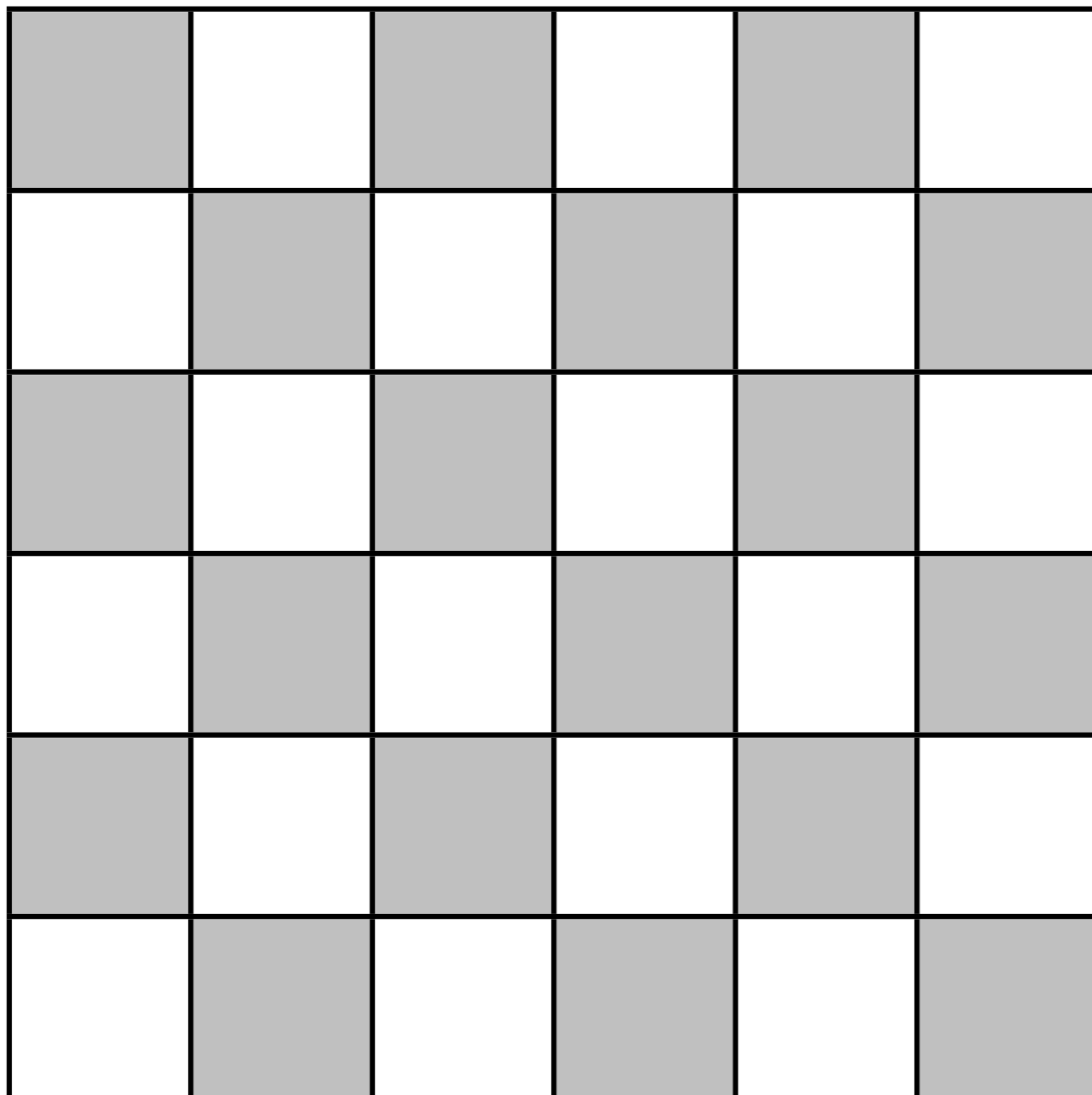
7

$$\begin{aligned} P(+\text{disease} \mid +\text{test1}, +\text{test2}) &= \\ P(+1, +2 \mid +d) P(+d) / P(+1, +2) &= \\ P(+12) = P(+12 \mid +d)P(+d) + P(+12 \mid -d)P(-d) &= \\ = 1 \cdot .9 \cdot .001 / (1 \cdot .9 \cdot .001 + .05 \cdot .01 \cdot .999) &= \\ = .643 & \end{aligned}$$

Test 1 seems better than test 2 -- why not use test 1 twice?

A: T1 is conditionally independent of T2 given disease state -- probably not true for T1 w/ itself

Independence



Geoff Gordon—10-701 Machine Learning—Fall 2013

8

for events: defined as $P(a \wedge b) = P(a)P(b)$, like rows/cols of checkerboard

$P(\sim a \wedge b) = P(b) - P(a \wedge b) = P(b) - P(a)P(b) = (1 - P(a))P(b) = P(\sim a)P(b)$

for r.v.s: $P(A, B) = P(A)P(B)$

shorthand: joint probability table = outer product of marginal probability tables

$P(A=a_i \wedge B=b_j) = P(A=a_i) P(B=b_j)$

intuition: knowing a or $\sim a$ tells us nothing about B

Bayes rule version: $P(A) = P(A|B)$

i.e., $P(A=a_i) = P(A=a_i | B=b_j)$ for all a_i, b_j w/ $P(B=b_j) > 0$

follows: $P(a|b) = P(a \wedge b)/P(b) = P(a)P(b)/P(b) = P(a)$, as long as $P(b) \neq 0$

some take this as definition of independence

Conditional independence



Geoff Gordon—10-701 Machine Learning—Fall 2013

9

conditional independence: $X \perp\!\!\!\perp Y \mid Z$ if, for all values of z_i , $P(X,Y|z_i)=P(X|z_i)P(Y|z_i)$
i.e., after we condition on z_i , X and Y are independent

$P(x, y \mid z) = P(x \mid z) P(y \mid z)$ for any values xyz of XYZ
this is a statement about a 3d table and two 2d tables

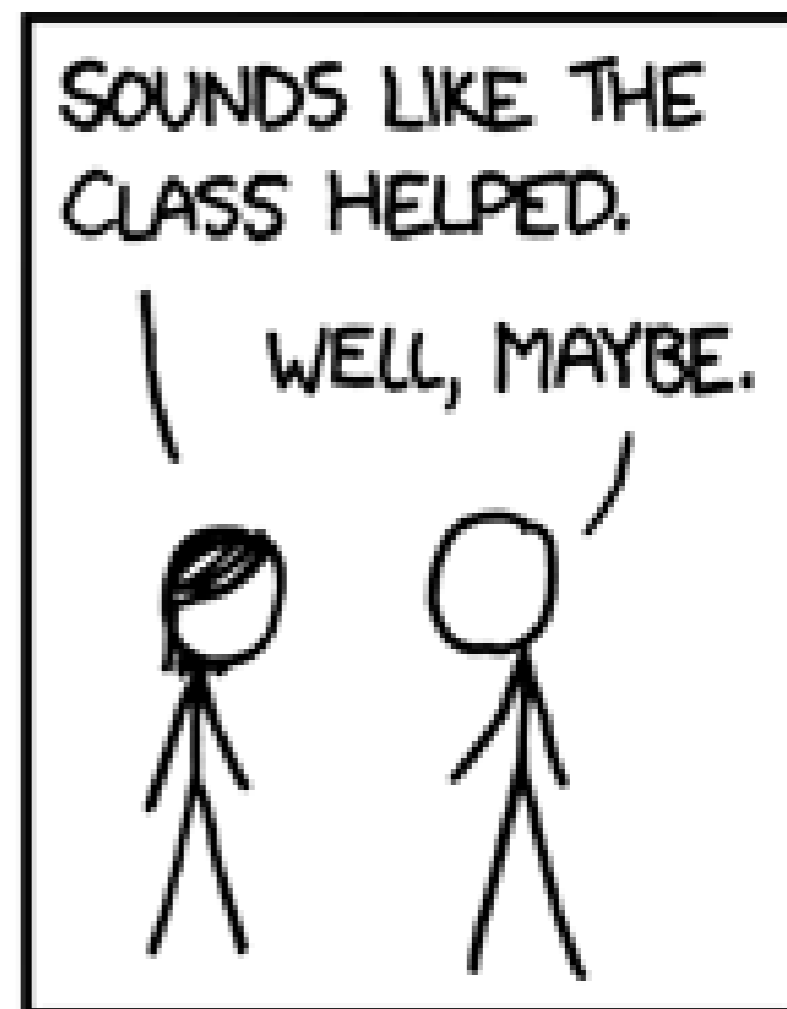
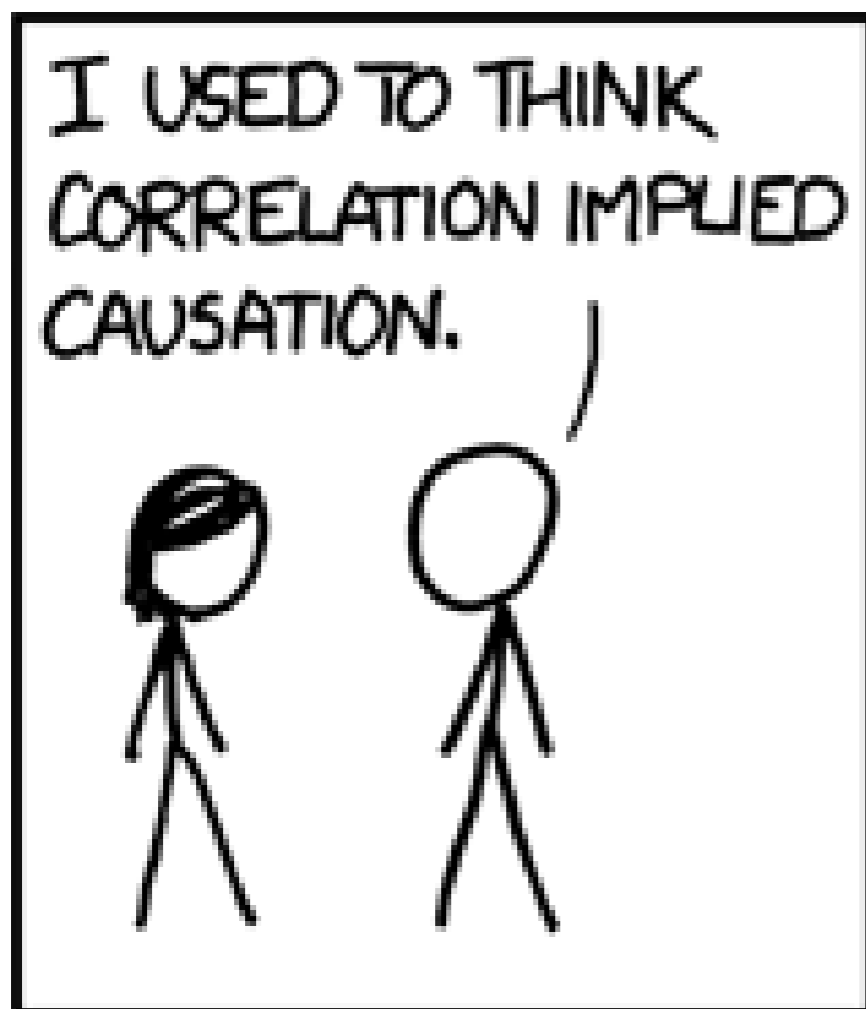
Conditionally Independent

London taxi drivers: A survey has pointed out a positive and significant correlation between the number of accidents and wearing coats. They concluded that coats could hinder movements of drivers and be the cause of accidents. A new law was prepared to prohibit drivers from wearing coats when driving.

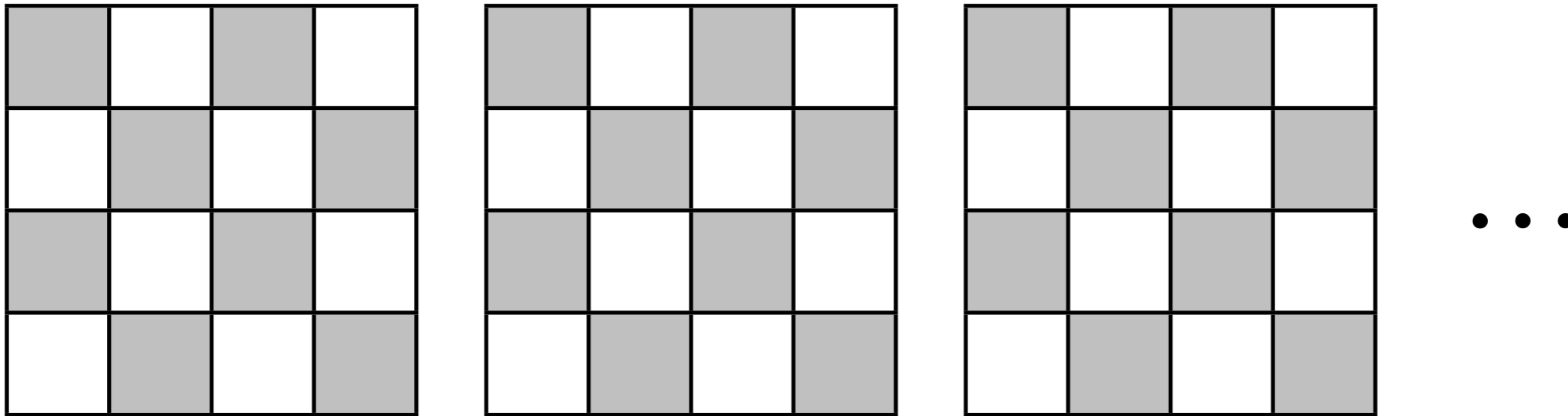
Finally another study pointed out that people wear coats when it rains...

slide credit: Barnabas

More on the importance of conditioning



Samples



Geoff Gordon—10-701 Machine Learning—Fall 2013

12

i.i.d. sample of an r.v.: N independent copies of the same checkerboard

sample is itself an r.v. in a bigger space (a $2N$ -dim hypercheckerboard)

atomic events are tuples of original atomic events

sample (r.v.) vs. population (the One True checkerboard)

statistic = any function of the sample

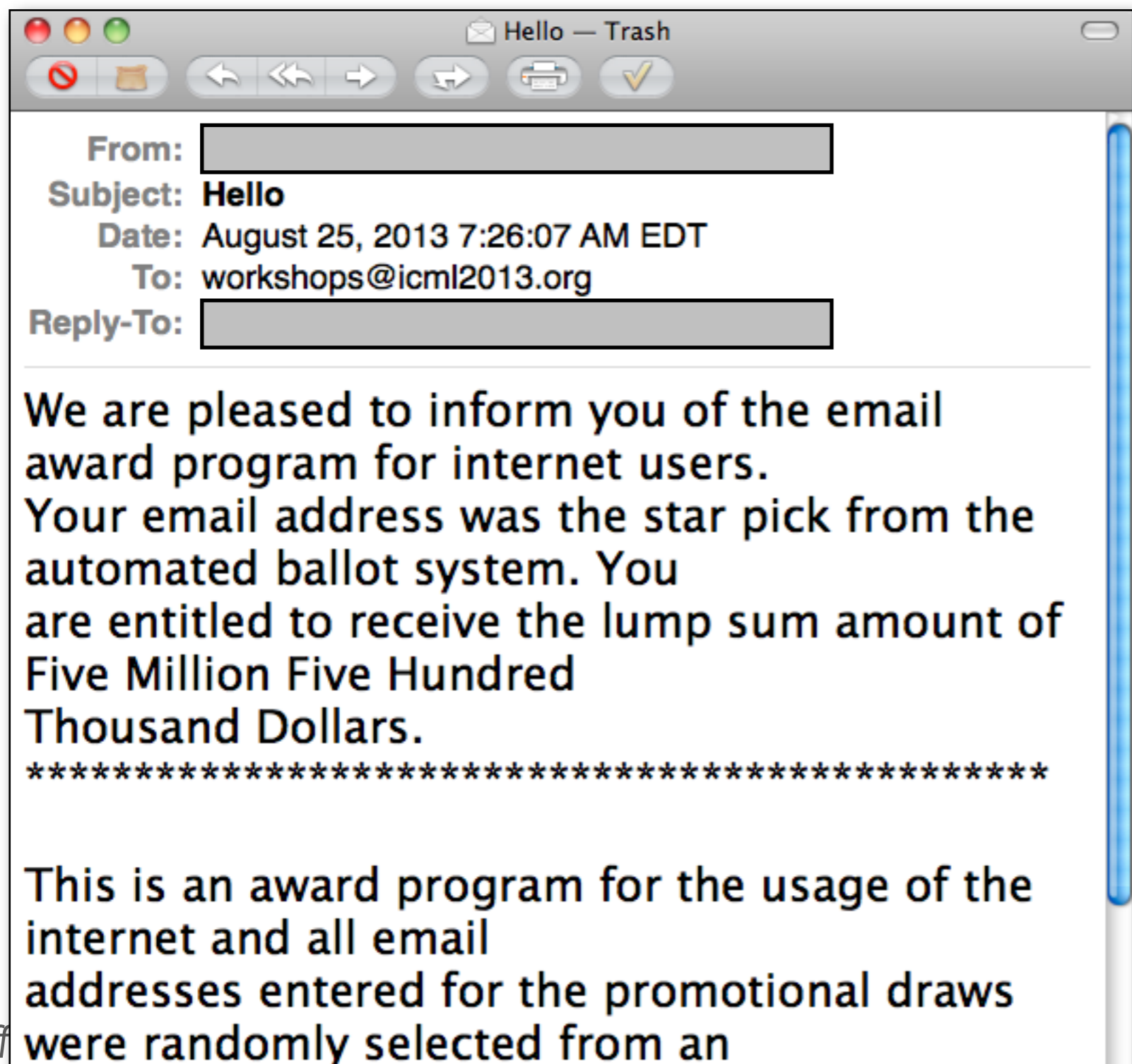
usu. order-independent (symmetric) but doesn't have to be

statistics are r.v.s

e.g., sample mean, variance

goal of statistics (big or small S): use statistics to find out something about the population

Recall: spam filtering



Geoff
classification problem: given data (x_i, y_i)
 N pairs, $x_i \in \{0,1\}^d$, $y_i \in \{0,1\}$ -- write $x_i = (x_{i1} \dots x_{id})$

produce rule which goes from feature $x \rightarrow$ predicted y
e.g.: spam filtering: x_{ij} = presence of word i in doc j

bag of words

A ridiculously naive assumption

- Assume:
- Clearly false:
- Given this assumption, use Bayes rule

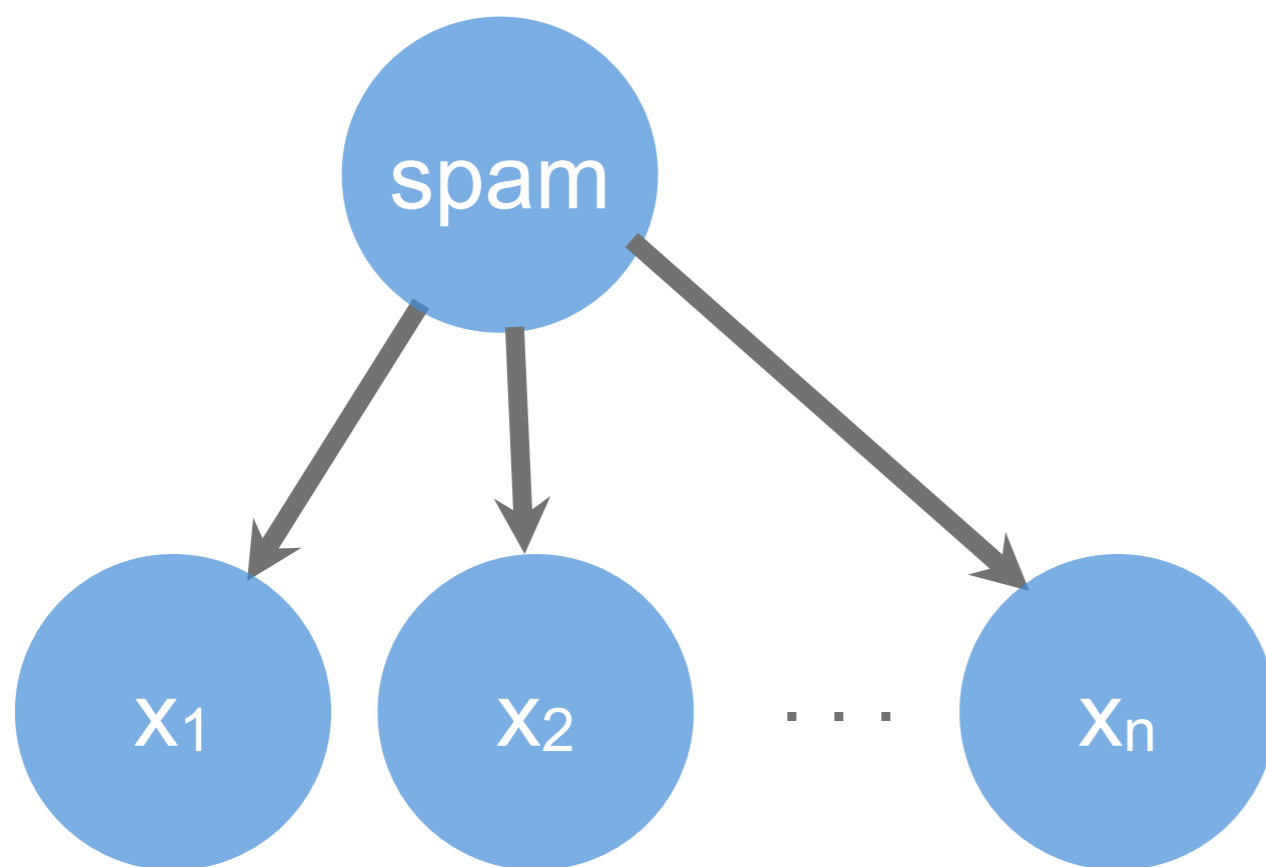
Geoff Gordon—10-701 Machine Learning—Fall 2013

assumption: $x_{ij} \perp x_{ik} \mid c_i$
for all $j, k \in 1..d, j \neq k$

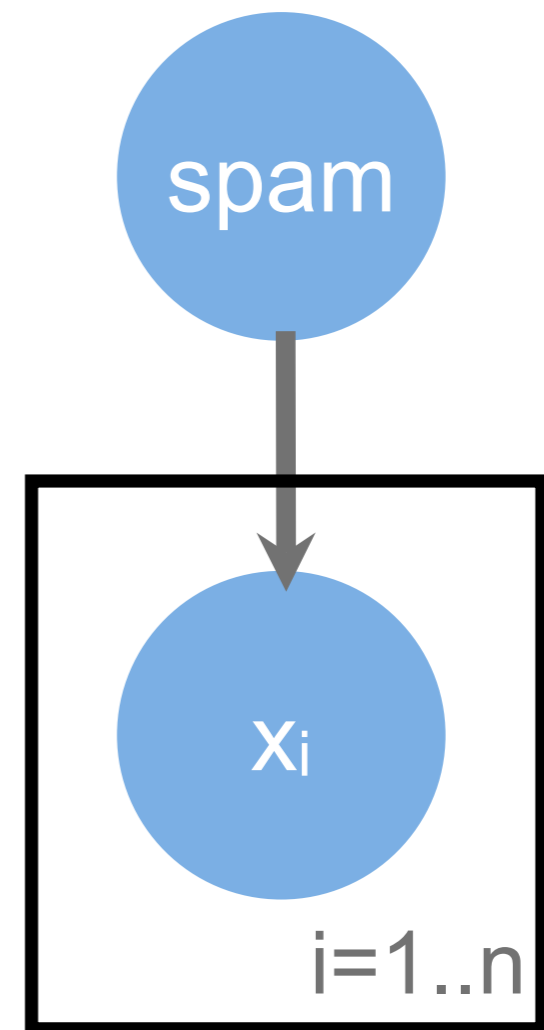
clearly false: "CMU" not independent of "Bayes"

given this "naive" assumption, use "Bayes" rule

Graphical model



$$P(X_1, \dots, X_d | \text{spam}) = \prod_{i=1}^d P(X_i | \text{spam})$$



Geoff Gordon—10-701 Machine Learning—Fall 2013

arrows spam → xi say xi depends on spam

lack of other arrows into xi: the xi's are conditionally independent given spam

“macro” or “for loop” shorthand: called a “plate model”

Naive Bayes

- $P(\text{spam} \mid \text{email} \wedge \text{award} \wedge \text{program} \wedge \text{for} \wedge \text{internet} \wedge \text{users} \wedge \text{lump} \wedge \text{sum} \wedge \text{of} \wedge \text{Five} \wedge \text{Million})$

Geoff Gordon—10-701 Machine Learning—Fall 2013

17

$$\begin{aligned} & P(\text{spam} \mid \text{email award program for internet users lump sum of Five Million}) \\ &= \frac{P(\text{email ... Million} \mid \text{spam}) P(\text{spam})}{[P(\text{email ... Million} \mid \text{spam}) P(\text{spam}) + P(\text{email ... Million} \mid \text{not-spam}) P(\sim\text{spam})]} \\ & \quad \text{(sum version of Bayes rule)} \\ &= \frac{P(\text{email} \mid \text{spam}) P(\text{award} \mid \text{spam}) \dots P(\text{Million} \mid \text{spam}) P(\text{spam})}{[P(\text{email} \mid \text{spam}) P(\text{award} \mid \text{spam}) \dots P(\text{Million} \mid \text{spam}) P(\text{spam}) + P(\text{email} \mid \sim\text{spam}) P(\text{award} \mid \sim\text{spam}) \dots P(\text{Million} \mid \sim\text{spam}) P(\sim\text{spam})]} \\ & \quad \text{(independence assumption)} \end{aligned}$$

suppose we know $P(\text{word } j \mid \text{spam})$ and $P(\text{word } j \mid \sim\text{spam})$ for all j
and suppose we know $P(\text{spam})$ and $P(\sim\text{spam})$
how? see slightly later

then above is easy to calculate!
now keep messages w/ $P(\text{spam}) < \text{threshold}$

adjust threshold based on user preference: chance of missing internet lottery win vs. wants to get work done

In log space

$$z_{\text{spam}} = \ln(P(\text{email} | \text{spam}) P(\text{award} | \text{spam}) \dots P(\text{Million} | \text{spam}) P(\text{spam}))$$

$$z_{\sim\text{spam}} = \ln(P(\text{email} | \sim\text{spam}) \dots P(\text{Million} | \sim\text{spam}) P(\sim\text{spam}))$$

Geoff Gordon—10-701 Machine Learning—Fall 2013

18

$$z_{\text{spam}} = \ln(P(\text{email} | \text{spam}) P(\text{award} | \text{spam}) \dots P(\text{Million} | \text{spam}) P(\text{spam}))$$

$$z_{\sim\text{spam}} = \ln(P(\text{email} | \sim\text{spam}) P(\text{award} | \sim\text{spam}) \dots P(\text{Million} | \sim\text{spam}) P(\sim\text{spam}))$$

result is then

$$P(\text{spam} | \dots) = \exp(z_{\text{spam}}) / [\exp(z_{\text{spam}}) + \exp(z_{\sim\text{spam}})]$$

$$= 1 / [1 + \exp(z_{\sim\text{spam}} - z_{\text{spam}})]$$

$$= 1 / [1 + \exp(-z)]$$

where $z = z_{\text{spam}} - z_{\sim\text{spam}}$

sigmoid or logistic function (sketch): $\text{logit}(z)$

big z : confident in spam

big $-z$: confident in $\sim\text{spam}$

Collect terms

$$Z_{\text{spam}} = \ln(P(\text{email} \mid \text{spam}) P(\text{award} \mid \text{spam}) \dots P(\text{Million} \mid \text{spam}) P(\text{spam}))$$

$$Z_{\sim\text{spam}} = \ln(P(\text{email} \mid \sim\text{spam}) \dots P(\text{Million} \mid \sim\text{spam}) P(\sim\text{spam}))$$

$$Z = Z_{\text{spam}} - Z_{\sim\text{spam}}$$

Geoff Gordon—10-701 Machine Learning—Fall 2013

19

= $\ln P(\text{spam}) - \ln(P(\sim\text{spam}))$ <-- call this b
+ $\ln P(\text{email} \mid \text{spam}) - \ln P(\text{email} \mid \sim\text{spam})$ <-- call this $w_{\{\text{email}\}}$
+ $\ln P(\text{award} \mid \text{spam}) - \ln P(\text{award} \mid \sim\text{spam})$ <-- $w_{\{\text{award}\}}$
+ ...

sum is over words in our message

i.e.: to classify message x_i , compute $b + \sum_{j \in \text{vocabulary}} w_j x_{ij}$

($x_{ij} = 0$ when word j not in message)

threshold on result: $b + \sum_j w_j x_{ij} \geq 0$

"linear discriminant"

"decision boundary" : $b + \sum_j w_j x_{ij} = 0$

Linear discriminant



Geoff Gordon—10-701 Machine Learning—Fall 2013

20

threshold on result: $b + \sum_j w_j x_{ij} \geq 0$

"linear discriminant"

"decision boundary" : $b + \sum_j w_j x_{ij} = 0$

Intuitions



Geoff Gordon—10-701 Machine Learning—Fall 2013

21

word of warning: nominally, $P(\text{spam})$ is $\text{logit}(z)$
really, $\text{logit}(z)$ usually close to 0 or 1 even in marginal cases
why? [failure of independence assumption]
fix: use $\text{logit}(\text{eps} * z)$ for appropriately-chosen eps (***)

what are highly useful/discriminative words?
 $\log(P(\text{word } j|\text{spam})/P(\text{word } j|\sim\text{spam}))$ far from zero
i.e., really low probability one class, moderate to high in the other
(since log is much more sensitive near 0)
but decent chance of actually seeing word j
suggested measure: $P(j)\log(P(j|s)/P(j|\sim s))$ -- looks like KL / relative entropy!
eg: "lottery" really low in legit emails -- not high in spam, but significantly nonzero

How to get probabilities?



Geoff Gordon—10-701 Machine Learning—Fall 2013

22

How to get the needed probabilities?

flip coin, count heads or tails

3H, 7T: estimate $P(H) = .3$

what if 0H, 2T?

hard zero probability seems bad

Hack: estimate $P(H) = (\#H + 0.5) / (\#H + \#T + 1)$

here, $0.5 / 3 = 1/6$: pretty sure but not hard zero

turns out not to be a hack after all...

called "Laplace smoothing", MAP for Dirichlet prior

for words, sample count: might expect it to be proportional to population prob
estimator: $\text{count}(\text{word}) / (\text{total} \# \text{ words})$

e.g., if we see "avocado" 10 times in 100k words, estimate $P(\text{avocado}) = 1/10k$

count 0: e.g., "syzygy" not in training set but appears in test doc
get $\log(0) = -\infty$ in answer, bad.

estimate $p(\text{word})$ as $(\lambda + \text{count}(\text{word})) / (\lambda * \text{words} + \text{total} \# \text{ words})$
small $\lambda > 0$

Improvements



n-grams, character string features, treat (multiple) headers vs. body differently, features of attachments, take account of body length, generalize among words using edit distance, stemming, POS-tagging, parsing, named entity resolution, collaborative spam-filters (if everyone's winning the lottery...), soft whitelists/blacklists (of domains, IP addresses, senders), ...

general rule: most important part of practical ML is to think of good sources of information, then figure out how to get algorithm to pay attention to them

Perceptron



Geoff Gordon—10-701 Machine Learning—Fall 2013

24

classification again:

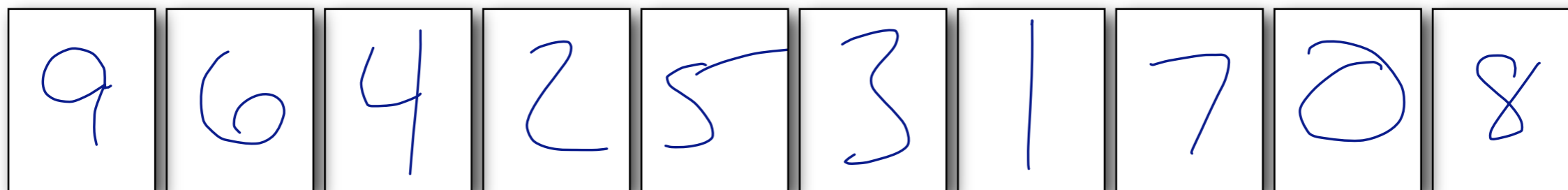
$X = ((x_1, y_1) \dots (x_N, y_N))$

$y_i \in \{-1, 1\}$ \leftarrow instead of 0/1 for convenience

assume $x_{ij} \in \{0, 1\}$ if desired, but works just as well w/ \mathbb{R}^d

ex: spam filtering, or digit classific'n

Digit recognition



digit recognition:

x_{ij} = pixel j intensity in image i

classify (say) 4 vs 9

Linear separability



Geoff Gordon—10-701 Machine Learning—Fall 2013

26

linear discriminant: $z_i = x_i \cdot w + b$
(note notation change: $x_i \cdot w = \sum_j x_{ij} w_j$; Alex wrote $\langle x_i, w \rangle$)

$z_i > 0$ when $y_i = +1$

$z_i < 0$ when $y_i = -1$

i.e., no errors and nothing on decision boundary

suppose data set X is separable; how can we find w in this case?

Simplify notation

Geoff Gordon—10-701 Machine Learning—Fall 2013

27

simplify notation: homogeneous rep:

$x_i \in \mathbb{R}^d$

assume last component $x_{id} = 1$ in all examples

if not, set $d := d+1$ and just append 1 to every example

now, add b to last component of w : say $w = w_{old} + [0 \ 0 \ 0 \ \dots \ b]$

$x_i \cdot w = x_i \cdot (w_{old} + [0 \ 0 \ 0 \ \dots \ b])$

$= x_i \cdot w_{old} + b = z_i$

benefit: don't need b in our notation any more

flip -ves:

write $u_i = y_i x_i \in \mathbb{R}^n$

$x_i \cdot w > 0$ when $y_i = +1$

$\implies u_i \cdot w > 0$

$x_i \cdot w < 0$ when $y_i = -1$

$\implies u_i \cdot w > 0$

so now, just write $u_i \cdot w > 0$ for all i

First algorithm: LP



Geoff Gordon—10-701 Machine Learning—Fall 2013

28

LP: we have a system of linear inequalities
almost an LP -- if it were, we'd be done (hand it to CPLEX and go for coffee)
but our ineqs are strict (LP software requires \geq not $>$)
suppose we know $\text{eps} = \min_i u_i \cdot w$
now we can write $u_i \cdot w \geq \text{eps}$ for all i
but we don't know eps -- trick: let $v = w/\text{eps}$
now $u_i \cdot v \geq 1$ for all i
solve for v -- and eps doesn't matter

remember this algorithm: we'll see another one like it when we get to support vector machines

#2: the “perceptron algorithm”

$$u_i \cdot w > 0 \quad \forall i$$

Geoff Gordon—10-701 Machine Learning—Fall 2013

29

simple algorithm:

start at $w = 0$

pick a random training point $u_i = y_i x_i$

if inequality is satisfied, great

if not, $w := w + u_i$

$u_i \cdot w_{\text{new}} = u_i \cdot w + u_i \cdot u_i > u_i \cdot w$

closer to satisfying ineq

repeat

This algorithm was implemented in 1958 on the “Mark I Perceptron”, which represented weights with variable resistors (i.e., dimmer switches), and used motors to turn the knobs for the update (!)

To analyze: draw a sphere around all u_i . (For convenience suppose radius 1; if not, renorm.)

Margin γ : $\max_{\|w\| \leq 1} \min_i u_i \cdot w$ (draw it)

Perceptron proof

Geoff Gordon—10-701 Machine Learning—Fall 2013

30

suppose w^* (unknown) separates w / margin γ , $\|w^*\|=1$

On every update,

$$\begin{aligned}w_{\text{new}} \cdot w^* &= (w + u_i) \cdot w^* \\ &= w \cdot w^* + u_i \cdot w^* \\ &\geq w \cdot w^* + \gamma\end{aligned}$$

So, after T mistakes, $w \cdot w^* \geq T \gamma$
(starts at 0, increases by γ each time)

$$\begin{aligned}w_{\text{new}} \cdot w_{\text{new}} &= (w + u_i) \cdot (w + u_i) \\ &= w \cdot w + 2 w \cdot u_i + u_i \cdot u_i \\ &\leq w \cdot w + 1\end{aligned}$$

So, after T mistakes, $w \cdot w \leq T$

$$\begin{aligned}T \gamma &\leq w \cdot w^* \\ &\leq \|w\| \|w^*\| \quad (\text{Cauchy-Schwartz}) \\ &= \sqrt{w \cdot w} \\ &\leq \sqrt{T}\end{aligned}$$

divide thru by \sqrt{T} and by γ :

$$\sqrt{T} \leq 1/\gamma$$