

Our map of data analytics

I have a question about my data

→ Compute a Statistic

How certain am I that this answer
is close to the population statistic?

→ Estimate uncertainty with
Bootstrapping (or classic
stats)

What if my data is not representative
enough of the population?

→ Fit a probability model

Single Variable questions

- "What is the average mass of a penguin?"
- "What is the probability C. Clark gets > 20 asts?"

"Joint" questions ("and")

- "What fraction of games does Caitlin Clark get > 10 Ast, > 10 pts ~~and~~ > 10 Reb?"
- "What is the mean bill len and flipper len?"

Conditional questions ("given")

- "What is the mean mass of Adelie penguins?"
- "What is probability of a penguin that weighs 5kg being a Gentoo?"

Our map of data analytics

I have a question about my data

→ Compute a Statistic

How certain am I that this answer
is close to the population statistic?

→ Estimate uncertainty with

Bootstrapping (or classic
stats)

What if my data is not representative
enough of the population?

→ Fit a probability model

Single Variable questions

"What is the average mass of a penguin?"

"What is the probability C. Clark gets > 20 assists?"

"Joint" questions ("and")

"What fraction of games does Caitlin Clark
get > 10 Ast, > 10 pts and > 10 Reb?"

"What is the mean bill len and flipper len?"

Conditional questions ("given")

"What is the mean mass of Adelie penguins?"

"What is probability of a penguin that
weights 5 kg being a Gentoo?"

Our map of data analytics

I have a question about my data

→ Compute a Statistic

How certain am I that this answer
is close to the population statistic?

→ Estimate uncertainty with

Bootstrapping (or classic stats)

What if my data is not representative
enough of the population?

→ Fit a probability model ↪

Single Variable questions

"What is the average mass of a penguin?"

"What is the probability C. Clark gets > 20 assists?"

"Joint" questions ("and")

"What fraction of games does Caitlin Clark
get > 10 Ast, ≥ 10 pts and > 10 Reb?"

"What is the mean bill len and tippyf len."

Conditional questions ("given")

"What is the mean mass of Adelie penguins?"

"What is probability of a penguin that
weights 5 kg being a Gentoo?"

Our map of data analytics

I have a question about my data

→ Compute a Statistic

How certain am I that this answer
is close to the population statistic?

→ Estimate uncertainty with

Bootstrapping (or classic stats)

What if my data is not representative
enough of the population?

→ Fit a probability model

Single Variable questions

- "What is the average mass of a penguin?"
- "What is the probability C. Clark gets > 20 assists?"

"Joint" questions ("and")

- "What fraction of games does Caitlin Clark get > 10 Ast, > 10 pts ~~and~~ > 10 Reb?"
- "What is the mean bill len and tippel len.?"

Conditional questions ("given")

- "What is the mean mass of Adelie penguins?"
- "What is probability of a penguin that weighs 5 kg being a Gentoo?"

Variable types and models

Type	Subtypes	Distributions
Nominal	Categorical General	Bernoulli Categorical ?
Ordinal	Ordinal	?
Quantitative	Integer Real	Positive Bounded → Binomial Positive Bounded → Exponential Bounded → Uniform Unbounded Normal Student's t

Probability Notation & terminology

Random Variable : \underline{X} → A value we haven't yet seen

e.g. The number of ast. Caitlin Clark gets in her next game

Support : What possible values could \underline{X} take?

Distribution : How probable is each value of \underline{X} ?

PDF/pmf : $p(x) = \text{Prob}[\underline{X} = x]$

function that defines these probabilities

"True distribution" : What we would get if we picked an observation at random from the population (unknown)

Model : An approximation to the true distribution that makes some simplifying assumptions (estimated)

Probability Notation & terminology

mean of a sample $\bar{X} = \sum_{i=1}^N \frac{1}{N} X_i$
Average of observations

mean of a distribution $E[\bar{X}] = \sum P(x) x$
[Expectation of a random variable] or
 $\int P(x) x dx$

Weighted average of all possible observations

Distributions w/ Multiple Variables

Data set w/ 2 variables PTS and Ast

Our next observation has 2 unknown values

X: number of PTS Y: number of Ast

Joint distribution: How probable is every combination?

Joint pdf / pmf: $p(x, y) = \text{Prob}[\underline{X}=x, \underline{Y}=y]$

$p(10, 12) \rightarrow$ "What is the probability of 10 pts and 12 Ast?"

Distributions w/ Multiple Variables

Marginal distribution: How probable is each value of \underline{X} ignoring \underline{Y}

Marginal Pdf/Pmf: $P(x) = \text{Prob}[\underline{X} = x]$

Independence: Condition that $P(x, y) = P(x)P(y)$
for all combinations

The values are unrelated

e.g. rolling 2 dice

Distributions w/ Multiple Variables

Conditional distribution: how probable is each value of X if we know Y ?

e.g. "If we know caitlin has 30 pts, we know she is playing well and may have more assists"

"If we know the penguin we're looking at is a gentoo, we know it's more likely to have a higher mass even before we weigh it"

Conditional pdf/pmf: $P(X|y) = \text{Prob}[X=x | Y=y]$

e.g. $P(5000 | \text{Gentoo})$ = What is the probability of our penguin being 5Kg if we know its Gentoo?

Probability Rules

$$P(X, Y) = P(X|Y) P(Y)$$

$$P(X) = \sum_y P(X, y) \quad (\text{discrete}) \quad \int P(x, y) dy \quad (\text{cont.})$$

$$P(X|Y) = \frac{P(X, Y)}{P(Y)}$$

$$P(Y|X) = \frac{P(X|Y) P(Y)}{P(X)}$$

Categorical distribution

If we want a distribution over a categorical variable, we need a prob. for each value

$$P(c) = p_c \quad \text{or} \quad P(x) = \sum_{c=1}^K I(x=c) p_c$$

$$\text{MLE: } p_c = \frac{1}{N} \sum_{i=1}^N I(x_i=c)$$

For 2 variables:

$$P(c, d) = p_{cd} \quad \text{or} \quad P(x, y) = \sum_{c=1}^k \sum_{d=1}^k I(x=c) I(y=d) p_{cd}$$

$$\text{MLE: } p_{cd} = \frac{1}{N} \sum_{i=1}^N I(x_i=c) I(y_i=d)$$

Bernoulli distribution

Categorical with \leq possible values $\{0, 1\}$

→ only need 1 probability

$$P(1) = p \quad P(0) = 1 - p$$

$$P(x) = p I(x=1) + (1-p) I(x=0) = p^x (1-p)^{1-x}$$

Conditional distribution

$$P(x|y) =$$

Compute probability for every subset

$$\text{MLE: } P_{c|d} = \frac{\sum_{i=1}^N I(x_i = c) I(y_i = d)}{\sum_{i=1}^N I(y_i = d)}$$

Subset w/
 $y = d$

In Pandas or

Not really any distinction from Sample Prop. (yet!)

Computed Variables

Our "variables" don't need to be columns of our original data frame
→ we can compute new ones

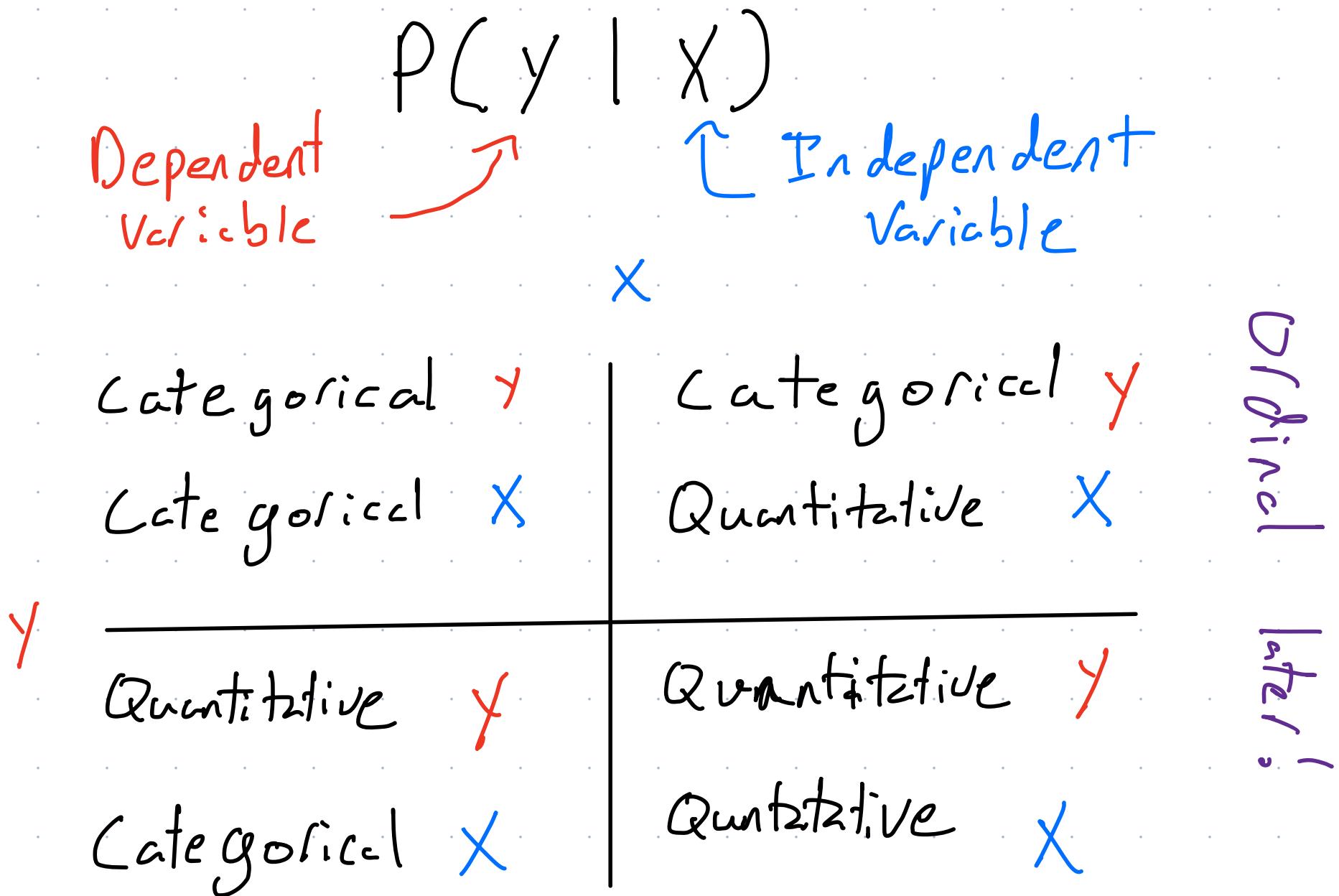
E.g. Quantitative Variable:

Caitlin Clark's assists

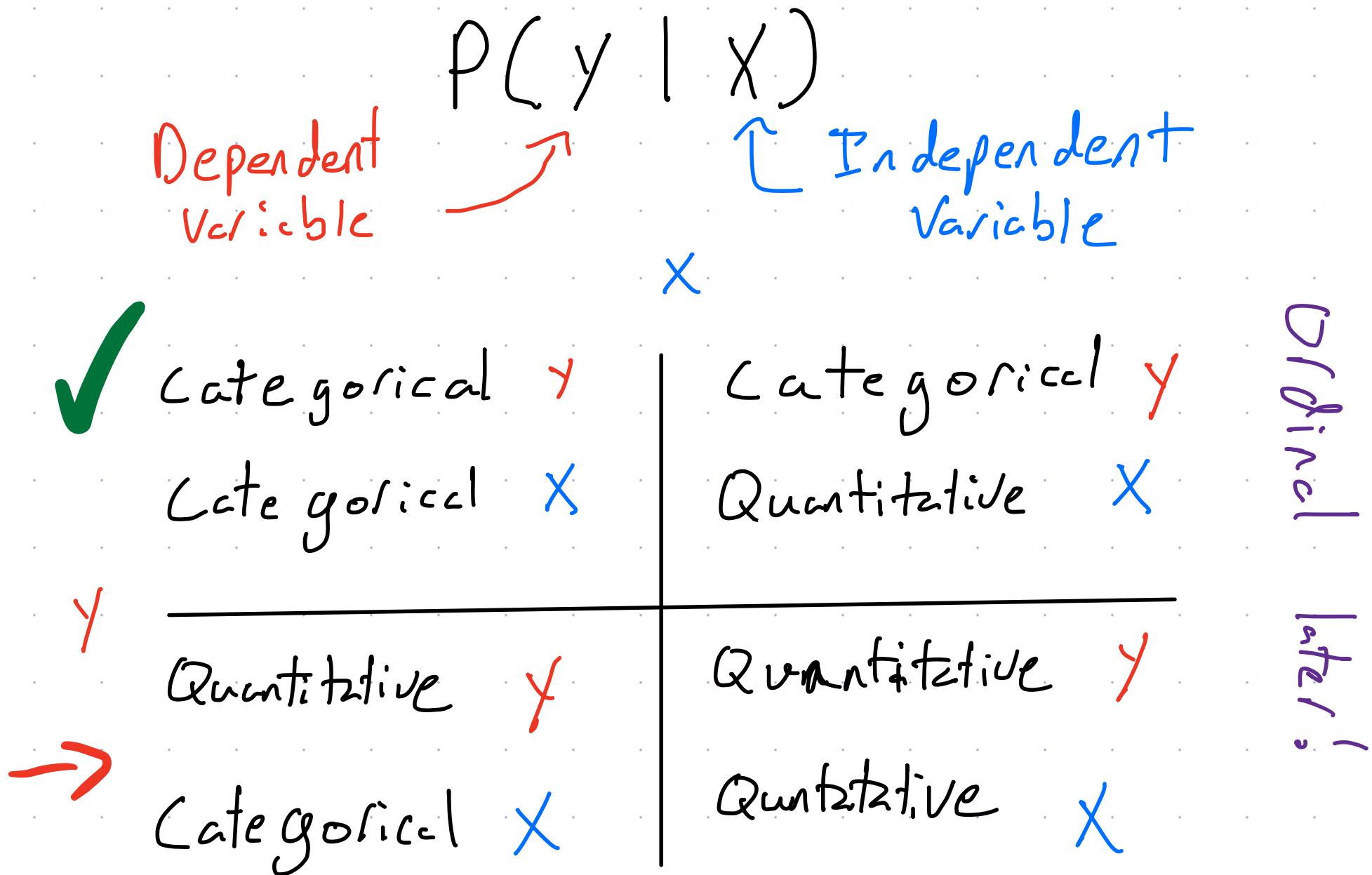
→ Categorical Variable:

Caitlin Clark's assists ≥ 10

Space of Conditional Models



Space of Conditional Models



P(Quantitative | Categorical)

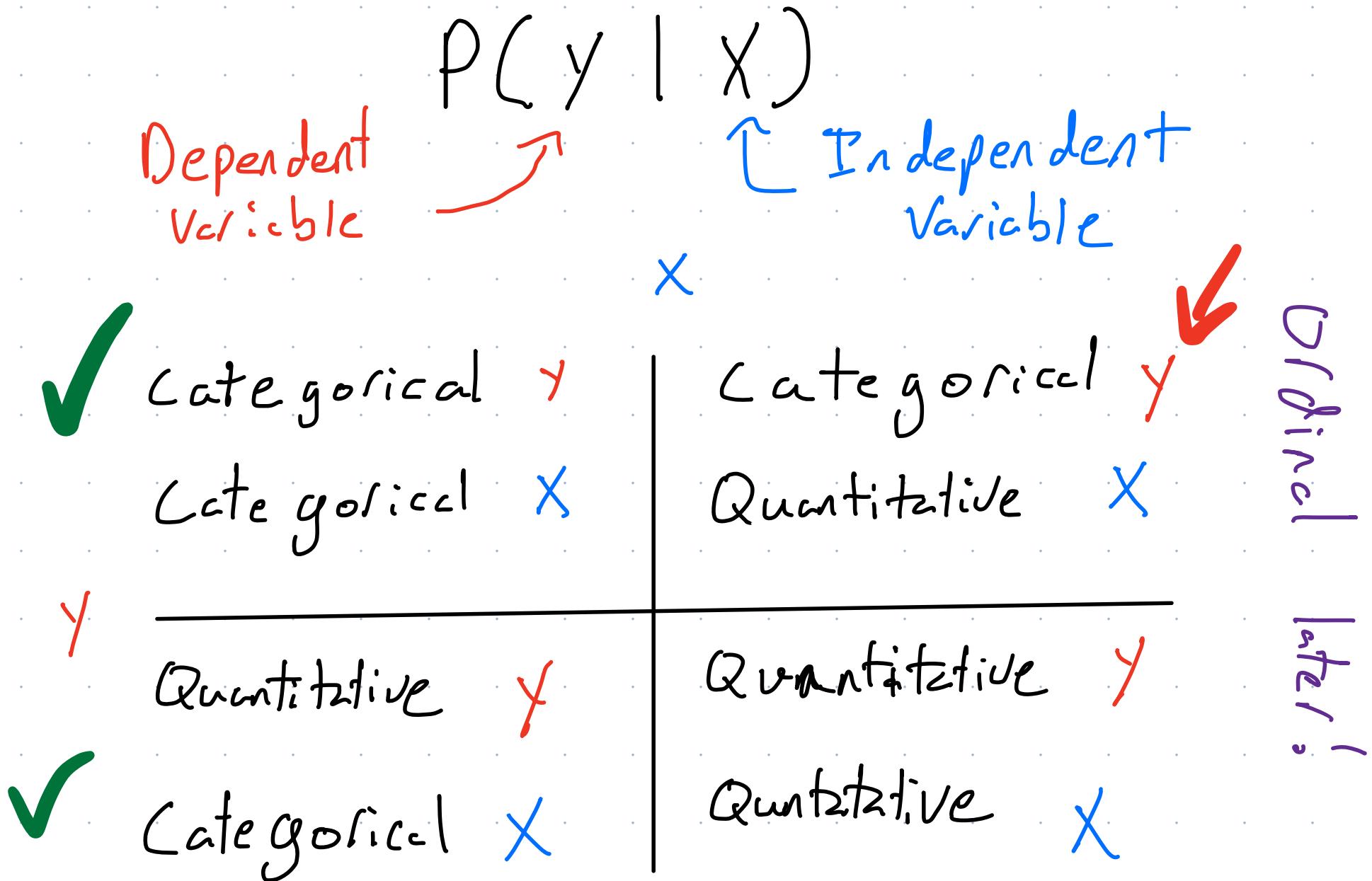
Straight forward!

- Choose appropriate distribution for quantitative variable
- Fit on each subset

E.g. Model Caitlin Clark's assists conditioned on Home Vs. Away

$$P(\text{Ast}|\text{home}) = \text{Poisson}(\text{Ast}; \hat{\lambda}), \hat{\lambda} = \frac{\sum_{i=1}^n I(x_i=\text{home}) y_i}{\sum_{i=1}^n I(x_i=\text{home})}$$

Space of Conditional Models



P(Categorical | Quantitative)

E.g. What's the probability the few /
won the game given how many points
Caitlin Clark scored?

"Compute the probability on every subset..."

$$P(\text{won}|0), P(\text{won}|1), P(\text{won}|2), \dots$$



Not easy! $P(\text{won}|20) = 1$

$$P(\text{won}|22) = 0 ?$$

P(Categorical | Quantitative)

Same problem that motivated models in the first place! \rightarrow We need to make more assumptions

Let's consider 2 possibilities:

1) Assume $P(x|y)$

is easy to model

$$\text{e.g. } P(\text{Pts}|\text{won}) = \text{Poisson}(\hat{\lambda}_w)$$

$$P(\text{pts}|\text{loss}) = \text{Poisson}(\hat{\lambda}_l)$$

2) Assume $P(y|x)$
changes predictably as x
changes

e.g.

$$P(y|x) = \frac{e^{ax+b}}{1 + e^{ax+b}}$$

$$P(\text{Categorical} \mid \text{Quantitative})$$

How does modeling $p(x|y)$ help us?

Bayes Rule!

Quantitative | Categorical



✓ $e^{-\zeta_j}$ to estimate!
Marginal

$$P(y|x) = \frac{P(x|y) P(y)}{P(x)}$$

$$P(x) = \sum_y P(x|y) P(y)$$

P(Categorical | Quantitative)

Ex.

$$P(\text{win} | \text{pts}) = \frac{P(\text{pts} | \text{win}) P(\text{win})}{P(\text{pts} | \text{win}) p(\text{win}) + P(\text{pts} | \text{loss}) p(\text{loss})}$$

$$P(\text{pts} | \text{win}) = \text{Poisson}(\hat{\lambda}_w), \quad \hat{\lambda}_w = \frac{\sum_{i=1}^N x_i I(y_i = \text{win})}{\sum_{i=1}^N I(y_i = \text{win})}$$

$x_i = \text{pts}, \quad y_i = \text{win}/\text{loss}$

$$P(\text{pts} | \text{loss}) = \text{Poisson}(\hat{\lambda}_l)$$

$$P(\text{win}) = \frac{\sum_{i=1}^N I(y_i = \text{win})}{N}$$

$$P(\text{loss}) = \frac{\sum_{i=1}^N I(y_i = \text{loss})}{N}$$

P(Categorical | Quantitative)

Wait! This is just a mixture model, But we know the assignments!

M-Step: Fit $\hat{P}(x|y)$, $P(y)$

Get probabilities for unknown y

E-Step: find $P(y|x)$

P(Categorical | Quantitative)

Assume $P(Y|X)$ changes predictably with X

Logistic regression → For 2 outcomes

Reminder: Let's define $P(Y=1|X=0) = b$
assume as X increases by 1, $P(Y|X)$ always
increases by α

so, $P(Y|X) = \alpha X + b$, But

$P(Y|X) \in [0, 1]$ or bad things happen ...

P(Categorical | Quantitative)

So, apply a function that maps $\mathbb{R} \rightarrow [0, 1]$

e.g. Sigmoid = $\sigma(x) = \frac{e^x}{1 + e^x}$

Logistic regression

$$P(y|x) = \frac{e^{ax+b}}{1 + e^{ax+b}}$$

← usually the
better choice

Probit regression

$$P(y|x) = \Phi(ax+b)$$

Normal cdf

P(Categorical | Quantitative)

How do we choose a and b ?

Maximum likelihood!

$$\underset{a, b}{\operatorname{argmax}} \prod_{i=1}^N P(Y_i | X_i) = \underset{a, b}{\operatorname{argmax}} \sum_{i=1}^N \log P(Y_i | X_i)$$

$$= \underset{a, b}{\operatorname{argmax}} \sum_{i=1}^N \log \frac{e^{ax+b}}{1 + e^{ax+b}}$$

No easy solution $\circlearrowleft \rightarrow$ Gradient descent

P(Categorical | Quantitative)

What Assumptions are we making for each?

Bayesian Model

$P(y|x)$ determined by
ratio of distance
to each mean

Not monotonic in x

Which is better? Depends on the data!

Can compute likelihoods

Logistic Regression

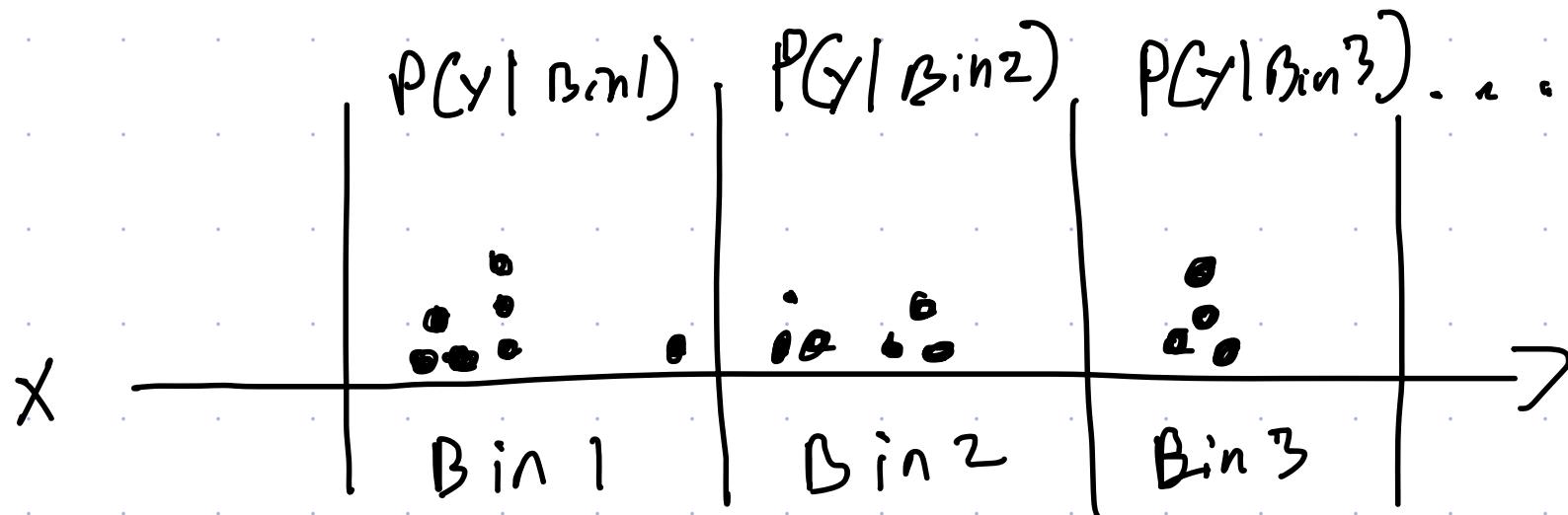
$P(y|x)$ determined by
magnitude of x

Monotonic in x

P(Categorical | Quantitative)

Other alternatives?

What about a histogram style approach?



→ Decision trees

(maybe) later in this course...

Space of Conditional Models

$$P(y \mid x)$$

Dependent Variable

I Independent Variable

Categorical Y

Cate gosiccl X

10

Categorical y

Quantitative X

Ordnung

Quantitative

Quantitative X

Categorical X

P(Quantitative | Quantitative)

Does our Bayesian method work?

No 

$$P(y|x) = \frac{P(x|y) P(y)}{P(x)}$$

Still would need to enumerate every y

$$P(\text{Quantitative} \mid \text{Quantitative})$$

Linear Regression:

$$P(y \mid x) \text{ is Normal } N(\mu, \sigma)$$

$$\text{Let mean } (\mu) \text{ of } P(y \mid x=0) = b$$

Assume that as x increases by 1

$$\mu \text{ increases by } a \rightarrow \mu = ax + b$$

[σ is constant]

$$P(y \mid x) = N(ax + b, \sigma^2) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{y - ax - b}{\sigma} \right)^2}$$

$$P(\text{Quantitative} \mid \text{Quantitative})$$

Picking a, b and σ

Maximum likelihood!

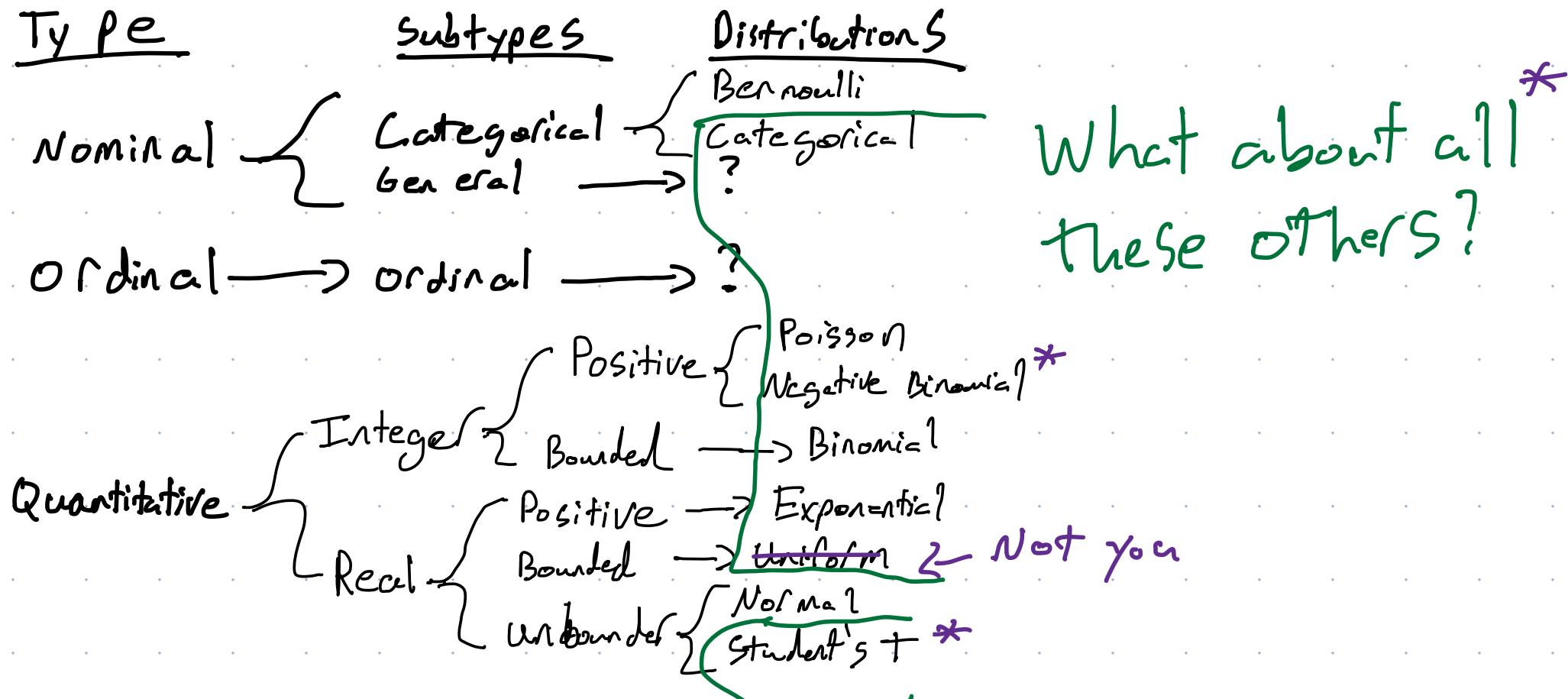
$$\underset{a, b}{\operatorname{argmax}} \prod_{i=1}^N P(y_i \mid x_i) = \underset{a, b}{\operatorname{argmax}} \sum_{i=1}^N \log P(y_i \mid x_i)$$

$$= \underset{a, b}{\operatorname{argmax}} \sum_{i=1}^N \log \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x_i - ax_i + b}{\sigma} \right)^2}$$

Here there is an "easy" solution

$$P(\text{Quantitative} \mid \text{Quantitative})$$

Linear regression lets us assume $p(y|x)$ is Normal. Logistic regression assumed Bernoulli



Generalized linear models

Mean of a Normal distribution $E[x] = \mu$

Linear regression $E[y|x] = \mu = ax + b$

Mean of a Bernoulli distribution $E[x] = p$

Logistic regression $E[y|x] = p = \sigma(ax + b)$

General Distribution

obey conditions for mean

Assume $E[y|x] = f(ax + b)$

Generalized linear models

What if we want $p(y|x)$ to be a Poisson?

Unconditional poisson $p(y) = \frac{\lambda^y e^{-\lambda}}{y!}$

mean: $E[y] = \lambda$

Poisson regression
Let $\lambda = ax + b$? No!
 $\lambda \geq 0$

$$\lambda = e^{ax+b}$$

$$f(x) = e^x$$

$$p(y|x) = \frac{(e^{ax+b})^y e^{-e^{ax+b}}}{y!}$$

Generalized linear models

Recipe:

Distribution: Normal, Bernoulli etc.

Structure: $\beta_1 x + \beta_0$

Link function: $g(\cdot)$

$$E[Y|X] = g^{-1}(\beta_1 x + \beta_0)$$

Categorical Regression

What if we want $p(Y|X)$ to be categorical w/ more than 2 outcomes?

e.g. $Y \in \{a, b, c\}$, no ordering

Categorical distribution

$$P(a) = P_a = .5 \quad P(b) = P_b = .3 \quad P(c) = P_c = .2$$

$$E[Y] = .5a + .3b + .2c = ?$$

"Dummy" Variables

Already saw we can compute new variables, compute 1 for each value!

$$Y \in \{a, b, c\}$$

$$\begin{aligned} Y \rightarrow & \quad A = I(Y=a) \in \{0, 1\} \\ & \quad B = I(Y=b) \in \{0, 1\} \\ & \quad C = I(Y=c) \in \{0, 1\} \end{aligned}$$

$$E[A] = P_a \quad E[B] = P_b \quad E[C] = P_c$$