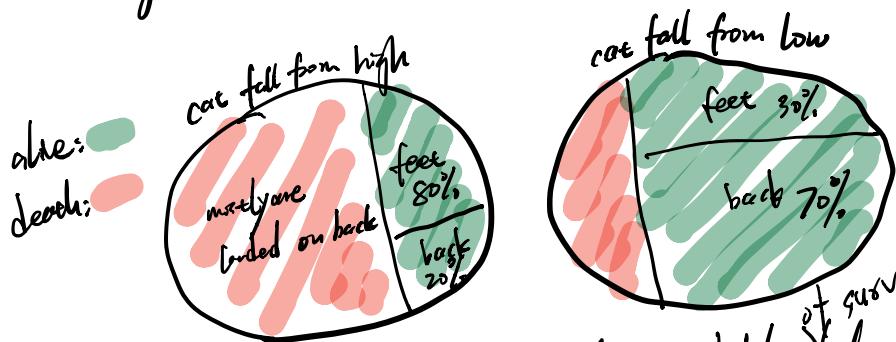


Cae Survival Experiment:

The longer the fall, the greater chance of survival.



car land on feet has a higher probability than car land on back even if they fall from higher place.

$$P(\text{survive} \mid \text{high. feet}) > P(\text{survive} \mid \text{low. back})$$

All red ones were sent to hospital.

- Alived ones were sent to hospital.
 - Among those ones fall from high. we say 70% of the cat land on feet survived so the $P(\text{survived} \mid \text{high}) = 70\% \times 80\% + 0\% \times 20\% = 56\%$
 $\Rightarrow 56\%$ of cat fall from high are survived in the hospital.
 - Among those ones fall from low. we say 80% of the cat land on feet survived, and 20% of cat land on back survived. so
 $P(\text{survived} \mid \text{low}) = 20\% \times 70\% + 80\% \times 30\% = 28\%$ in the hospital.
 - The reason is because most ones fall from high are dead. If they are alive, they most likely land on feet. so they will have a high survival rate. However, most ones fall from low are not dead and sent to hospital. Lots of them land on back. so they will have a low survival rate after being sent to hospital. so the general probability

of survival of cats fall from low is lower than those fall from high.
 However, she thinks if most cats fall from high are clearly dead,
 so they were not sent to hospital.

Reformulate the problem:

$$P(\text{Survive 1 week} \mid \text{fell})$$

Questions to ask:

1. Does the vet data represent all cat who fell?

2. If not, is the data missing at random? (MAR)

MAR in this case: $P(\text{go to vet} \mid \text{fell}) = P(\text{go to vet} \mid \text{fell, state of cat})$

We can use the law of total probability to work this out:

$$P(\text{Survive 1 week} \mid \text{fell}) = P(\text{Survive 1 week, ! clearly dead} \mid \text{fell}) \\ + P(\text{Survive 1 week, clearly dead} \mid \text{fell})$$

$$P(A, B \mid C) = P(A \mid C) P(B \mid A, C)$$

$$P(\text{Survive 1 week, ! clearly dead} \mid \text{fell}) = P(! \text{ clearly dead} \mid \text{fell}) \\ \times P(\text{Survive 1 week, ! clearly dead} \mid \text{fell})$$

$$\underbrace{P(\text{Survive 1 week, clearly dead} \mid \text{fell})}_{\theta} = P(\text{clearly dead} \mid \text{fell}) \\ \times \underbrace{P(\text{Survive 1 week} \mid \text{clearly dead, fell})}_{\theta}$$

$$\therefore \boxed{P(\text{Survive 1 week} \mid \text{fell}) = P(! \text{ clearly dead} \mid \text{fell}) \times P(\text{Survive 1 week} \mid ! \text{ clearly dead, fell})}$$

Equation 1

Assume the sampling mechanism: If ! clearly dead \Rightarrow go to net thus.

$$P(\text{Survive 1 week} \mid \text{fell, sampled}) = P(\text{Survive 1 week} \mid \text{fell, ! clearly dead})$$

So the reporting implies: $P(\text{Survive 1 week} \mid \text{fell, sampled}) = P(\text{Survive 1 week} \mid \text{fell})$

But in reality: $P(\text{survive 1 week} \mid \text{fell, sampled})$ Based on equation ①

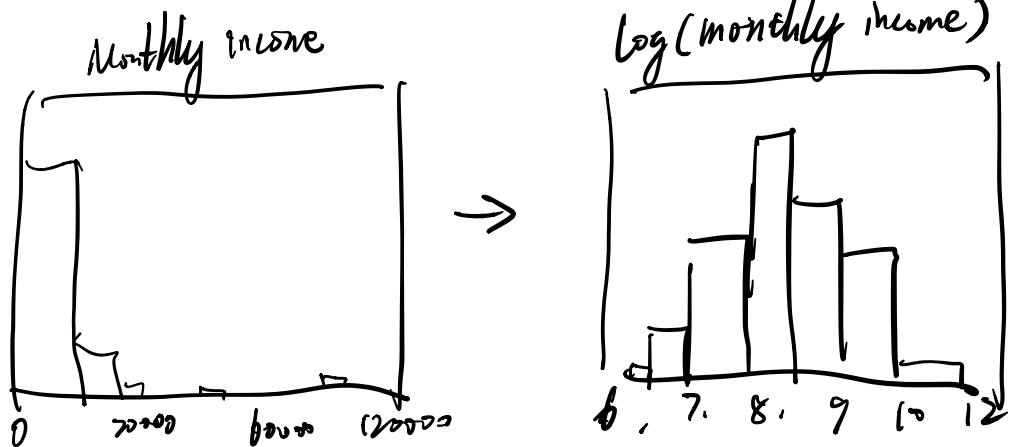
$$= P(\text{survived 1 week} \mid \text{fell}) / P(\text{! clearly dead} \mid \text{fell})$$

The degree of bias depends on how low $P(\text{! clearly dead} \mid \text{fell})$ is.

If fell from high, $P(\text{! clearly dead} \mid \text{fell})$ will be low. so the $P(\text{survive 1 week} \mid \text{fell, sampled})$ will be high. which explains the observation of the reporting

Knowing your data:

Understand its shape. The histogram is a great tool for this:



Outliers:

It's important to distinguish between an outlier from a normal process of the data-generating mechanism versus an outlier caused by some sort of failure within the system.

Missing data:

First to know whether the data is missing at random.
look at the following page.

Why distributions are important when training ML model?

Some ML models are designed to work best under distribution assumptions. Therefore, knowing which distribution we are working with, can help us to identify which models are best to use.

Numerical data → Discrete (number of students) → PMF
continuous (heights and weights) → PDF

- PMF gives the probability that a value equal to certain value.
- PDF itself is not probability, they need to be integrated over the given range to get the probability.

Data Cleaning:

HANDLING OUTLIERS/NULLS

What should one do when faced with these issues?
The easiest answer in data science is *it depends...*

If missing/bad at random...

- If the occurrence is rare, delete observations (most extreme)
- Can also impute/replace with average or median for that feature.

Otherwise...

- Impute with some constant (usually the average or median), create a dummy variable to indicate missing/bad
- Exploit multi-collinearity, i.e., use a model to estimate $E[\text{Missing Val}|X]$.

1. ~~remove~~ → delete

2. Impute avg or median

3. Create missing column which could be indicative

4. Estimate the missing value using other variable (Exploit multi-collinearity?)

No matter your preferred technique, this becomes a testable design choice. i.e,

1. Identify multiple competing methods for imputation and/or outlier cleaning
2. Train a model for each method (on the same training data)
3. Evaluate each method against the same out-of-sample validation data
4. Choose the best performing

Copyright Brian d'Alessio. All rights reserved

1. Identify multiple competing methods for imputation
 2. Train a model for each method using the same data
 3. evaluate each method's performance against the same one-of-sample data
 4. choose the best performing method.
-

selection bias: the sample can not represent the general population.

selection bias affects the generalizability of results and potentially the identifiability of model parameters.

Generalizability:

Does your model represent the population of interest?
more loosely: average s.d. $y-f(x)$

ex: estimate the avg(height) of the world population, but you went to elementary school, so the avg(school population) can not represent the whole population.

Identifiability: Any estimate derived from the sample that do not contain our population of interest is not identifiable.

ex: you want to estimate the avg(height) of women, but your sample contains all men, so the avg(sample) cannot represent the population of interest (no women) this is our population of interest.

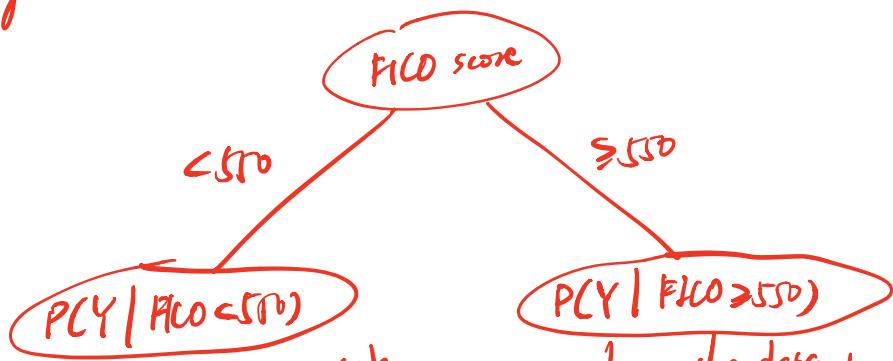
Credit Risk Modeling:

Goal: Predict PC Default in 6 mos (application data)

Method:

1. Sample new credit card users, log app data
2. Observe for 6 months, log if user defaults
3. Build a predictive model on sampled observations.

In this case, since we don't approve the credit card for user with $FICO < 550$, we don't have the data for these users. so our sample only contains users with > 550 and have missing data.



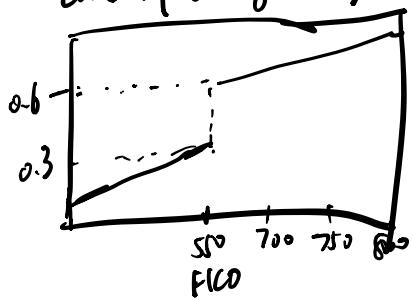
Because our data does not contain this data, so it's not identifiable.

If the model is truly linear, we can extrapolate to this region, but there could be an expensive assumption.

Solution to the problem:

1. we can accept everybody, and the result will be ideal for learning, but it could be a huge investment, which could break the bank

2. take a hybrid approach. balance the need to learn and the cost of losing money.



could have high variance for the estimate of user ($\text{FICO} < 550$)

Concept drift (Non-stationary)

sample used for model may be different from the population over time.
Example causes: seasonality, economic cycles, holiday promotions

Should we build a new model?

If the distribution has changed we should probably rebuild.

What data to use?

having more doesn't mean having the right data.

T1 (TV campaign)	T2 (no campaign)	T3 (no campaign)
------------------	------------------	------------------

unseen

we know T3 would more like T2

we can pull a small portion of T2 data as the **evaluation set**,
and build two models with ^{and} and without T1 to compare the performance.
the rest of T2

1. Monitor predictive measure constantly.
2. Retrain as often as possible (if the model degrade)
3. Test balance between data recency and data volume.