# Maximum Likelihood Estimation w/ missing data

MLE is a method to estimate the parameters of a probabilistic model given data by finding the _likeliest_ parameter values for those observations given the model.

Suppose we have a set of $n$ data (observations) $X_1, X_2, \ldots, X_n$. (iid *)

Our model is $Pr(X_1, X_2, \ldots, X_n | \Theta)$, the joint probability density given a set of parameters $\Theta$.

→ since iid, this reduces to

$$Pr(X_1, \ldots, X_n | \Theta) = Pr(X_1 | \Theta) Pr(X_2 | \Theta) \cdots Pr(X_n | \Theta) = \prod_{i=1}^{n} Pr(X_i | \Theta)$$

→ This lets us say how probable the $\{X_i\}$ are, given $Pr(\cdot)$ and $\Theta$.

→ We can _reverse_ this to ask: how _likely_ is $\Theta$, given the $Pr(\cdot)$ and the $\{X_i\}$.

$$\Rightarrow \mathcal{L}(\Theta; X_1, \ldots, X_n) = \prod_{i=1}^{n} Pr(X_i | \Theta)$$

"variable" of $\mathcal{L}$

"parameters" of $\mathcal{L}$

often one works w/ the "log-likelihood"
$$\ell = \ln \mathcal{L} = \sum_{i=1}^{n} \ln Pr(X_i | \Theta)$$

Now we can "plug in" different $\Theta$'s into $\mathcal{L}$ to see which is best. (actually optimize $\Theta$ to maximize $\mathcal{L}$ either analytically or numerically)

Example: binary data $X_i = 0$ or $X_i = 1$
model $Pr(X_i = 1 | \Theta) = \Theta \leftarrow \Theta \equiv$ prob of a "1".

Likelihood: $\mathcal{L}(\Theta | \{X_i\}) = \prod_{i=1}^{n} \Theta^{X_i} (1-\Theta)^{1-X_i}$    $\Theta_{MLE} = \dfrac{\sum_{i=1}^{n} X_i}{n}$

Missing Data    ML is fairly convenient at handling missing values.

Suppose we have _two_ variables $X_i, Y_i$ for each observation. Then

$$\mathcal{L}(\Theta) = \prod_{i=1}^{n} Pr(X_i, Y_i | \Theta)$$    as before.

Now suppose $X_i$ is missing for some observations:

→ First $m$ observations have both $X_i$ and $Y_i$

→ Remaining $n-m$ observations have $Y_i$ only.

For an $i$ w/ missing data, the probability for $Y_i$ only can be computed from its <u>marginal</u> distribution:

$$f(Y_i | \theta) = Pr(\cdot, Y_i | \theta) = \sum_{X} Pr(x, Y_i | \theta) \quad \text{if } x \text{ discrete}$$

— sum over all possible values of $x$. (integrate if $x$ is continuous)

Likelihood of $\theta$ given entire sample is now:

$$\mathcal{L}(\theta) = \prod_{i=1}^{m} Pr(x_i, Y_i | \theta) \prod_{i=m+1}^{n} f(Y_i | \theta)$$

$\underbrace{\qquad\qquad}$ complete observations

$\underbrace{\qquad\qquad}$ incomplete observations

ML strengths: effective for MCAR, MAR! Software available for linear models.

weakness: Requires a model for joint distribution of all variables w/ missing data. Often, assume multivariate normal distribution, but this is not feasible for all problems.

# Imputation

→ Fill in the blanks → Replace the missing observations w/ (hopefully) reasonable estimates.

## Marginal Mean Imputation

→ Replace any missing values for variable $X_j$ w/ the mean value of the non-missing $X_j$ entries.

## Conditional Mean Imputation:

If only one variable $X_j$ has missing values, learn a model for it from the other variables:

$$X_j = g(X_1, X_2, \ldots, X_{j-1}, X_{j+1}, \ldots, X_p)$$

and replace the missing entries $x_{ij}$ with $g(x_{i1}, x_{i2}, \ldots x_{i,j-1}, x_{i,j+1}, \ldots x_p)$

→ Gets really complicated when multiple $X_j$'s have missing values → which is almost always true.

## Last observation Carried Forward

• specific to repeated observations taken over time (longitudinal studies)

| observation time | | | |
|---|---|---|---|
| unit | $t_1$ | $t_2$ | $t_3 \cdots$ |
| 1 | 1.7 | 2.0 | ? ? ? |
| 2 | 4.1 | 3.7 | 4.0 $\cdots$ |
| 3 | | | |

$\Rightarrow$

| observation time | | | |
|---|---|---|---|
| unit | $t_1$ | $t_2$ | $t_3 \cdots$ |
| 1 | 1.7 | 2.0 | 2.0 2.0 $\cdots$ |
| 2 | | | |
| 3 | | | |

## Logical Rules and Relatedness

Suppose yearly income is missing from a survey, but, sometimes a related question — # months unemployed — is given. If 12 months unemployed, reasonable to impute yearly income = 0.

# Simple random imputation

For each missing value for variable $j$, replace it with a randomly chosen value from the observations where $j$ is observed. Sample w/ replacement if multiple observations are missing variable $j$.

→ OK if few observations are missing (but then you can just use listwise deletion) but this uses none of the information available regarding how $p$ relates to other variables

# Random (regression) imputation.   (RI)

Perform conditional mean imputation (typically w/ a linear model) but then simulate uncertainty in each imputed measurement by drawing a random noise value $\epsilon$ and adding it onto the imputed value.

$$X_{ij} = g\left(X_{i1}, X_{i2}, \ldots, X_{i(j-1)}, X_{i(j+1)}, \ldots X_p\right) + \epsilon$$

where typically $\epsilon \sim N(0, \sigma^2)$ and $\sigma^2$ is ~~appropriately~~ estimated from the same procedure used to fit $g$.

# Multiple Imputation

Next step up from random imputation.

→ want to reflect uncertainty both in sampling variation (like Random Imputation) but also uncertainty in $g$ itself.

<u>Idea:</u> take $M$ different models $g_1, g_2, \ldots, g_m$, perform R.I. on each, to make $M$ different complete datasets $X_1, \ldots, X_M$.   (typically $M = 5 - 10$)

→ complete observations are duplicated across $X$'s, missing values introduce a <u>new source</u> of variability.

Then estimate $\hat{\beta}_i$ on each, and pool these together in some way.

## Pooling:

$$Y = \alpha_i + \beta_i X_i, \quad i = 1, M$$

Suppose each $\hat{\beta}_i$ is a simple linear regression. Pooling here means combining the estimated regression coefficients

$$\hat{\beta} = \frac{1}{M} \sum_{i=1}^{M} \hat{\beta}_m$$

Where this becomes tricky is we also have uncertainties ⟶ standard errors associated with each $\hat{\beta}$ and these need to be pooled as well. (variation both within and between imputations is needed)