

Data Quality and Patient Matching Metrics

| METRICS / KEY PERFORMANCE INDICATORS FOR MATCHING | | | | | |
|---|---|-------------|--------------------------------|---|---|
| HIMSS developed a set of key performance indicators (KPIs) that allow an organization to evaluate its patient matching processes and technology and make continuous improvements. | | | | | |
| Metric | Definition | Calculation | | Numerator | Denominator |
| EMPI Database Activity Rate (EDAR) | This rate provides the relative activity of the database. It provides the ratio of the total number of EMPI algorithm matching evaluations (TEM) performed in a given timeframe to the overall size of the EMPI database (EDS). | EDAR = | Total EMPI Matches (TEM) x 100 | Total EMPI Matches (TEM): The total number of potential pair candidates presented by the EMPI algorithm for a given period of time. A computer generated duplicate report will identify the number of candidate pairs for matches based on the algorithm rules or thresholds. | EMPI Database Size (EDS): The total number of records or lines stored in the database. This includes individual |
| | | | EMPI Database Size (EDS) | | |

categories such as unique individuals, duplicates, persons, non-eligible individuals for security access control, quality control, names, pseudonyms

| | | | | | |
|--|--|--------|---|--|--|
| | | | | | s, test names, unrecorded records, and abandoned records, registrations, among others. |
| EMPID at database size Duplicate Rate (EDDR) | This is the percent of paired records in the database that are potential duplicates or multiples. This measure reports a duplicate rate prior to research and validation of the records paired by the algorithm. It is commonly referred to as the Database Duplicate Rate, Duplicate Percentage or sometimes Pair Rate. | EDDR = | <div>EMPI Database Duplicates (EDD) x 100</div> <div>EMPI Database Size (EDS)</div> | EMPI Database Duplicates (EDD): The numeric count of records that are potential duplicates within the database. The EDD is calculated by subtracting the EPP (unduplicated person count in the database) from the EDS (total lines in database). | EMPID at database size (EDS): The total number of records |

or
li
n
e
s
st
or
e
d
in
th
e
d
a
t
a
b
a
s
e.
T
h
i
s
in
cl
ud
e
s
in
di
vi
d
u
al
c
at
e
g
or
ie
s
s
u
ch
a
s
u
ni
qu
e
in
di
vi
d
u
al
s,
d
u
pl
ic
at
e
s,
p
er
s
o
n
n
el
is
t
in
g
s
fo
r
s
e
c
ur
it
y
a
c

| | | | | | |
|--------|---|-------|---|---|---|
| | | | | | c e s s c o n t r o l, q u a l i t y c o n t r o l n a m e s, p s e u d o n a m e s, t e s t n a m e s, u n- r e c o n c i l e d r e c o r d s, a n d a b a n d o n e d r e g i s t r a t i o n s, a m o n g o t h e r s. |
| D u | This is the ratio of newly created duplicate records (numerator) to the | DCR = | Total number of individual duplicate patient records x 100 | Total EMPI Matches (TEM): The total number of potential pair candidates presented by the EMPI | T ot |

| | | | | |
|--|---|--|--|--|
| pl ic at e C r e at io n R at e (D C R) | <p>opportunity to create a duplicate through various encounters with patients (denominator) in a given period of time. Opportunities to create duplicates would include scheduling events, registration, preregistration, office visits, among others. Denominators may vary by organization depending on the registration and scheduling solutions employed. The total number of individual duplicate patient records is obtained by dividing the Total EMPI Matches (TEM) by 2.</p> | <p>Total Registrations Performed (TRP)</p> | <p>algorithm for a given period of time. A computer generated duplicate report will identify the number of candidate pairs for matches based on the algorithm rules or thresholds.</p> | al R e gi st ra ti o n s P er fo r m ed (T R P): T h e to ta l n u m b er of re gi st ra ti o n a ct iv iti es (r e gi st ra ti o n, pr er e gi st ra ti o n, s c h e d ul in g, of fi c e vi si ts , a m o n g ot her |
|--|---|--|--|--|

s) performed in a given time period. These activities may vary by organization at independence on the registration and scheduling including solutions and examples

| | | | | | |
|--|--|-------|---|---|---|
| | | | | | y e d a s w e l l a s t h e i r b u s i n e s s m o d e l. |
| T r u e M a t c h R a t e (T M R) | True Match Rate (TMR): The True Match Rate is the ratio of the number of true match pairs (TMP), as determined after manual validation, to the total number of EMPI matching evaluations (TEM) presented by the algorithm plus those identified by other business processes. The formula for computing the TMR is the number of True Match Pairs (TMP) divided by the total number of potential pair candidates (TEM). This figure provides information on the effectiveness of the algorithm in making matches. A low true match rate may indicate that the programmed matching thresholds may need to be fine-tuned or adjusted. | TMR = | <div>True Matched Pairs (TMP) x 100</div> <div>Total EMPI Matches (TEM)</div> | True Matched Pairs (TMP): The number of pairs generated by the algorithm as well as external business processes that are found, after manual validation, to be confirmed as matched pairs. Sometimes referred to as Adjusted Matched Pairs (AMP). | T o t a l E M P I M a t c h e s (T E M): T h e t o t a l n u m b e r o f p o t e n t i a l p a i r c a n d i d a t e s p r e s e n t e d b y t h e E M P I a l |

gorithm for a given period of time. A computer generated at added duplicate report will identify the number of candidates pairs for match based on the al

| | | | | | |
|---|--|--------|---|---|---|
| | | | | | or it h m ru le s or th re s h ol d s. |
| F al s e P o si ti v e M a t c h R a t e (F M R) | False Positive Match Rate (FPMR): Sometimes referred to as False Match Rate or False Positive Rate. The percentage of incorrectly matched candidate pairs over a given period. This measures the percentage of invalid pairs that have been incorrectly paired by the algorithm. The False Match Rate is the incidence of False Matches made by the algorithm that have been confirmed or validated as not being the same individual. It is computed by dividing the number of False Positive Match Pairs (FPMP) by the total number of potential pair candidates (TEM). This figure provides information on the effectiveness of the algorithm in making matches. A high false match rate may indicate that the programmed matching thresholds are too permissive in their matching criteria and may need to be fine-tuned or adjusted. | FPMR = | False Positive Matched Pairs (FPMP) x 100 <hr/> Total EMPI Matches (TEM) | False Positive Matched Pairs (FPMP): The number of candidate pairs generated by the algorithm that are found, after manual validation, not to be matched pairs. These are sometimes referred to as False Positives or False Matched Pairs. It is the number of incorrect matched pair determinations made by the algorithm in a given period of time. | T o t a l E M P I M a t c h e s (T E M): T h e t o t a l n u m b e r o f p o t e n t i a l p a i r c a n d i d a t e s p r e s e n t e d b y th e E M P I a l g o r i t h m f o r a g |

ven period of time. A computer generated duplicated report will identify the number of candidates paired for match these based on the algorithm rules or the

| | | | | | |
|---|---|-------|---|--|---|
| | | | | | re s h o l d s. |
| F a l s e N e g a t i v e (N o n - M a t c h) R a t e (F N R) | False Negative (NonMatch) Rate (FNR): This reports the percent of incorrect EMPI Non-match decisions made in a given time frame. It is the percentage of candidate pairs who should have been matched but were not. These pairs were discovered during the course of business over a given period of time and were not identified by the algorithm. This measure reflects the matching status or decision after review and validation. This is a manual calculation. It must be pointed out that the FNR result will most probably reflect an incidence of unmatched records much lower than what actually exists in the database. This provides a view of the algorithm effectiveness in discovering matched pairs. Algorithm tuning may be required to reduce the incidence of unmatched pairs. | FNR = | False Negative 'non-match' Pairs (FNMP) x 100 Total EMPI Matches (TEM) | False Negative 'non-match' Pairs (FNMP): The number of incorrect EMPI non-match decisions made in a given period of time. The number of candidate pairs identified by the algorithm thresholds to be non-matches but after manual validation are determined to be matched pairs. These are most often found through normal business processes such as patient, physician, or scheduler report. These are not identified by the algorithm. In other words, this is the number of incorrect EMPI non-match decisions made. This will result in an otherwise unidentified duplicate remaining in the database. This is sometimes called a false negative count or pair. | T o t a l E M P I N o n - m a t c h e s (T E N M): T o t a l n u m b e r o f t r u e m a t c h e s m i s s e d b y t h e a l g o r i t h m a n d i d e n t i f i e d d u r i n g n o r m a l o p e r a t i o n p r o c |

essess that require additional validation of records. This is a measure of algorithm performance. Note: There is no easy way to capture this but this is a

| | | | | | |
|---|-------|---|--|---|---|
| | | | | | key metric that reflects the accuracy of the algorithm performance. |
| Indeterminate Match Rate (IMR): This is also known as ambiguous match rate. Of the total number of evaluations performed by the algorithm (TEM), the percent that were found to be indeterminate matches after validation. These are matches where the pair of candidate records offered by the algorithm did not have sufficient information to make a clear determination of whether or not they were the same individual. This may indicate a number of different factors such as data quality or data capture challenges, business process variation, "old" data that was not adequately managed, or simply business factors that influenced where an organization sets a matching threshold. | IMR = | <div>Indeterminate Match Pair (IMP) x 100</div> <div>Total EMPI Matches (TEM)</div> | Indeterminate Matched Pairs (IMP): After manual review, the number of algorithm candidate pairs whose identities could not be validated as being the same individual are called an indeterminate match decision. | Total EMPI Matches (TEM): The total number of potential pairs of candidates | |

s
p
r
e
s
e
n
t
e
d
b
y
t
h
e
E
M
P
I
a
l
g
o
r
i
t
h
m
f
o
r
a
g
i
v
e
n
p
e
r
i
o
d
o
f
t
i
m
e
. A
c
o
m
p
u
t
e
r
g
e
n
e
r
a
t
e
d
d
u
p
l
i
c
a
t
e
r
e
p
o
r
t
w
i
l
l
i
d
e
n
t
i
f
y
t
h
e
n
u
m
b
e
r
o
f
c
a
n
d
i
d
a
t
e
p
a
i
r
s
f
o
r

| | | | | | |
|---|--|-------|--|---|---|
| | | | | | m a t c h e s b a s e d o n t h e a l g o r i t h m r u l e s o r t h e s h o l d s. |
| M a t c h i n g A c c u r a c y R a t e (M A R) | Matching Accuracy Rate (MAR): This is the overall accuracy rate of the demographic matching process over a given period of time. | MAR = | <div>Total Match Pair (TMP) + Total NonMatch Pair (TNMP) x 100</div> <div>Total EMPI Matches (TEM)</div> | True Matched Pairs (TMP): The number of pairs generated by the algorithm as well as external business processes that are found, after manual validation, to be confirmed as matched pairs. Sometimes referred to as Adjusted Matched Pairs (AMP). | T o t a l E M P I M a t c h e s (T E M): T h e t o t a l n u m b e r o f p o t e n t i a l p a i r c a n d i d a t e s p r e s e n t e d b |

y
th
e
E
M
P
l
al
g
or
it
h
m
fo
r
a
gi
ve
n
p
er
io
d
of
ti
m
e.
A
c
o
m
p
ut
er
ge
ne
rat
ed
d
u
pl
ic
at
e
re
p
or
t
w
ill
id
en
t
if
y
th
e
n
u
m
be
r
of
c
a
n
di
d
at
e
p
ai
rs
fo
r
m
at
ch
es
ba
s

| | | | | | |
|---------------------------|---|-------|---|---|---|
| | | | | | ed on the algorithm rules or thresholds. |
| Matching Error Rate (MER) | This is the overall error rate of the demographic matching process over a given period of time. | MER = | <div>False Positive Match Pair (FPMP) + False Positive NonMatch Pair (FNMP) x 100</div> <div>Total EMPI Matches (TEM)</div> | False Positive Matched Pairs (FPMP): The number of candidate pairs generated by the algorithm that are found, after manual validation, not to be matched pairs. These are sometimes referred to as False Positives or False Matched Pairs. It is the number of incorrect matched pair determinations made by the algorithm in a given period of time. | Total EMPI Matches (TEM): The total number of potential pair candidates presented by the EMPI algorithm |

or it h m fo r a gi v e n p e r i o d of t i m e. A c o m p u t e r g e n e r a t e d d u p l i c a t e r e p o r t w i l l i d e n t i f y t h e n u m b e r of c a n d i d a t e p a i r s f o r m a t c h e s b a s e d o n t h e a l g o r

| | | | | | it h m ru le s or th re s h ol d s. |
|---|--|-------------|----------------------|--------------------------------|--|
| Other metrics/performance indicators | | | | | |
| M e t r i c | Definition | Calculation | | Numerator | D e n o m i n a t o r |
| P r e c i s i o n / P V a k a C l a s s i f i c a t i o n A c c u r a c y | <p>The proportion of true matches that were found out of the total matches found. Precision is the number of correct results divided by the number of all returned results.</p> <p>In other words, it is the number of correct positive results divided by the number of positive results predicted by the classifier.</p> | Precision = | $\frac{TP}{TP + FP}$ | Number of true positives found | N u m b e r o f t o t a l m a t c h e s f o u n d |
| R e c a l l a k a S e n s i t i v i t y | <p>How many from the matches found are real matches or the percent of all relevant documents that is returned by the search.</p> <p>In other words, it is the number of correct positive results divided by the number of all relevant samples (all samples that should have been identified as positive).</p> | Recall = | $\frac{TP}{TP + FN}$ | Number of true positives found | N u m b e r o f p o t e n t i a l t r u e m a t c h e s |

| | | | | | |
|------------------------|---|---|---|--|--------------------|
| F1 Score | <p>F1 Score is the harmonic mean between precision and recall. The range for F1 Score is [0, 1]. It tells you how precise your classifier is (how many instances it classifies correctly), as well as how robust it is (it does not miss a significant number of instances).</p> <p>High precision but lower recall, gives you an extremely accurate, but it then misses a large number of instances that are difficult to classify. The greater the F1 Score, the better is the performance of our model. Mathematically, it can be expressed as :</p> | $F = 2 *$ | <div>Precision * Recall</div> <hr/> <div>Precision + Recall</div> | Precision * Recall | Precision + Recall |
| | The probability of a correct classification. | $(1 - \text{Error}) =$ $\text{Pr}(C) =$ | <div>(TP + TN)</div> <hr/> <div>(PP + NP)</div> | | |
| Specificity | The ability of the test to correctly rule out the disease in a disease-free population. | Specificity = | <div>TN</div> <hr/> <div>(TN + FP)</div> | | |
| Logarithmic Loss | <p>Logarithmic Loss or Log Loss, works by penalising false positives. It works well for multi-class classification. When working with Log Loss, the classifier must assign probability to each class for all the samples.</p> <p>Suppose, there are N samples belonging to M classes, then the Log Loss is calculated:</p> | $\text{Logarithmic Loss} = \frac{-1}{N} \sum_{i=1}^N \sum_{j=1}^M y_{ij} * \log(p_{ij})$ | | <p>where, y_{ij}, indicates whether sample i belongs to class j or not p_{ij}, indicates the probability of sample i belonging to class j</p> <p>Log Loss has no upper bound and it exists on the range [0,). Log Loss nearer to 0 indicates higher accuracy, whereas if the Log Loss is away from 0 then it indicates lower accuracy. In general, minimizing Log Loss gives greater accuracy for the classifier.</p> | |
| Area Under Curve (AUC) | <p>Area Under Curve(AUC) is one of the most widely used metrics for evaluation. It is used for binary classification problem. AUC of a classifier is equal to the probability that the classifier will rank a randomly chosen positive example higher than a randomly chosen negative example. Before defining AUC, let us understand two basic terms :</p> <p>True Positive Rate (Sensitivity) : True Positive Rate is defined as TP / (FN+TP). True Positive Rate corresponds to the proportion of positive data points that are correctly considered as positive, with respect to all positive data points.</p> <p>False Positive Rate (Specificity) : False Positive Rate is defined as FP / (FP+TN). False Positive Rate corresponds to the proportion of negative data points that are mistakenly considered as positive, with respect to all negative data points.</p> <p>False Positive Rate and True Positive Rate both have values in the range [0, 1]. FPR and TPR both are computed at threshold values such as (0.00, 0.02, 0.04,, 1.00) and a graph is drawn. AUC is the area under the curve of plot False Positive Rate vs True Positive Rate at different points in [0, 1].</p> <p>As evident, AUC has a range of [0, 1]. The greater the value, the better is the performance of our model.</p> | <p style="text-align: center;">Receiver operating characteristic example</p> <p style="text-align: right;">ROC curve (area = 0.79)</p> | | | |

| | | | | |
|-------------------------|--|---|--|--|
| Mean Absolute Error | <p>Mean Absolute Error is the average of the difference between the Original Values and the Predicted Values. It gives us the measure of how far the predictions were from the actual output. However, they don't give us any idea of the direction of the error i.e. whether we are under predicting the data or over predicting the data. Mathematically, it is represented as :</p> | $MeanAbsoluteError = \frac{1}{N} \sum_{j=1}^N y_j - \hat{y}_j $ | | |
| Mean Squared Error(MSE) | <p>Mean Squared Error(MSE) is quite similar to Mean Absolute Error, the only difference being that MSE takes the average of the square of the difference between the original values and the predicted values. The advantage of MSE being that it is easier to compute the gradient, whereas Mean Absolute Error requires complicated linear programming tools to compute the gradient. As, we take square of the error, the effect of larger errors become more pronounced than smaller error, hence the model can now focus more on the larger errors.</p> | $MeanSquaredError = \frac{1}{N} \sum_{j=1}^N (y_j - \hat{y}_j)^2$ | | |