

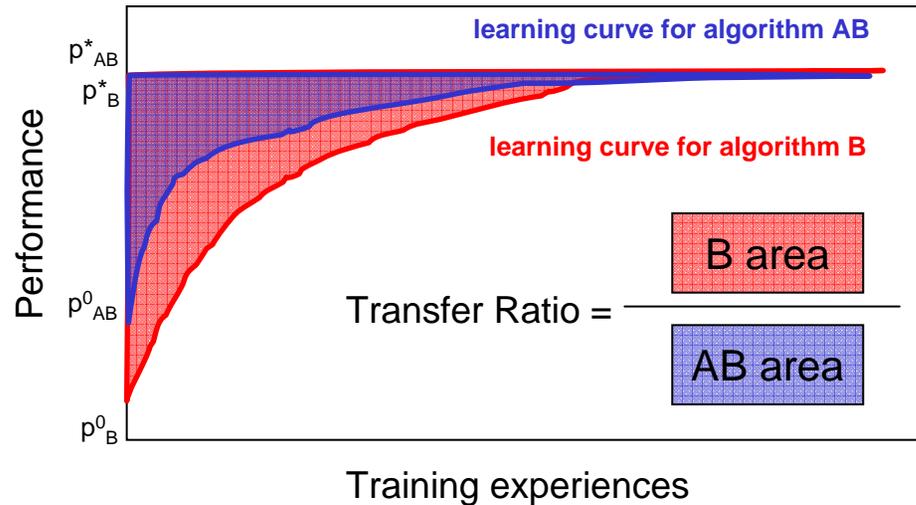


Proposed Metrics for Transfer Learning

Version 4
April 30, 2006

Tom Dietterich, editor

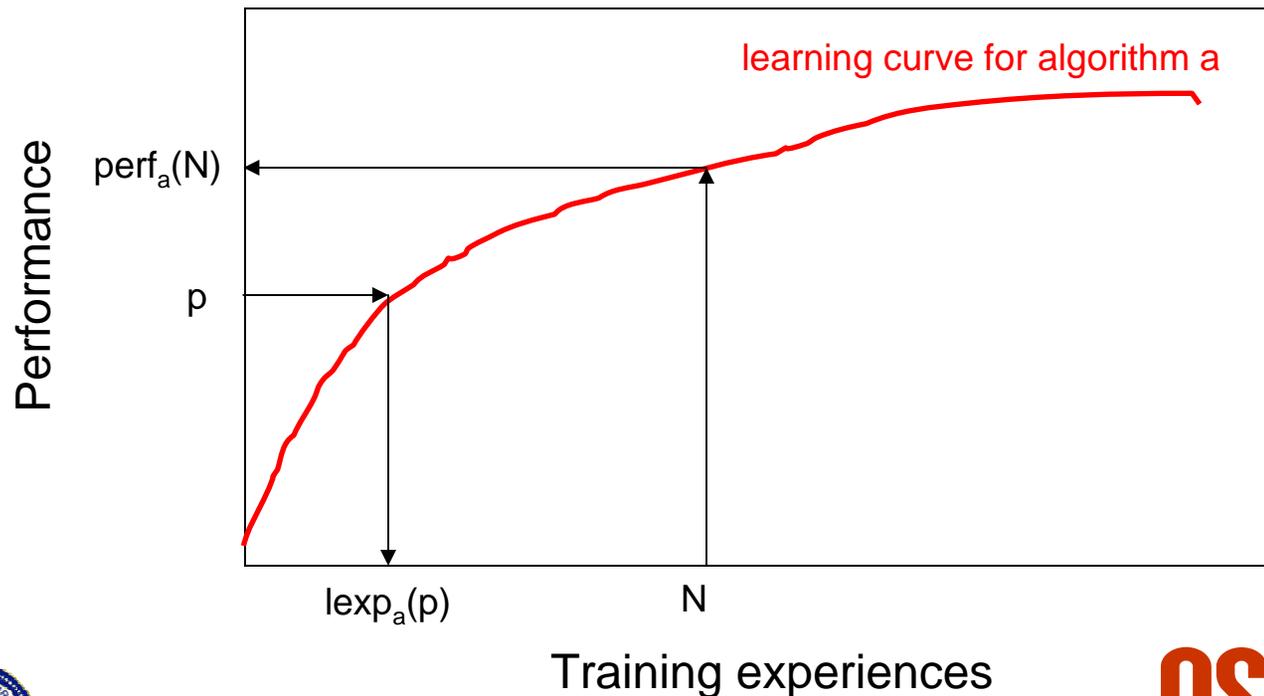




- Hard to explain: no well-defined units
- Large dependence on correctly estimating the asymptotic performance p_{AB}^*
- Infinite when $p_{AB}^* > p_B^*$
- Does not give an instantaneous value at each sample size. This can hide many interesting phenomena
- Alternate minimum TR metric only measures the worst-case. Can hide excellent performance at small-to-medium sample sizes

- Average speedup
 - advantage: easy to interpret
 - disadvantage: can be infinite even after removing obvious cases
- Average relative reduction in amount of training data required
 - advantage: also easy to interpret
 - advantage: avoids most infinite answers

- Let $p = \text{perf}_{\text{alg}}(N)$ be the performance of algorithm “alg” after N training “experiences”
- Let $N = \text{lexp}_{\text{alg}}(p)$ be the inverse function: the number of training experiences N required to achieve target level p



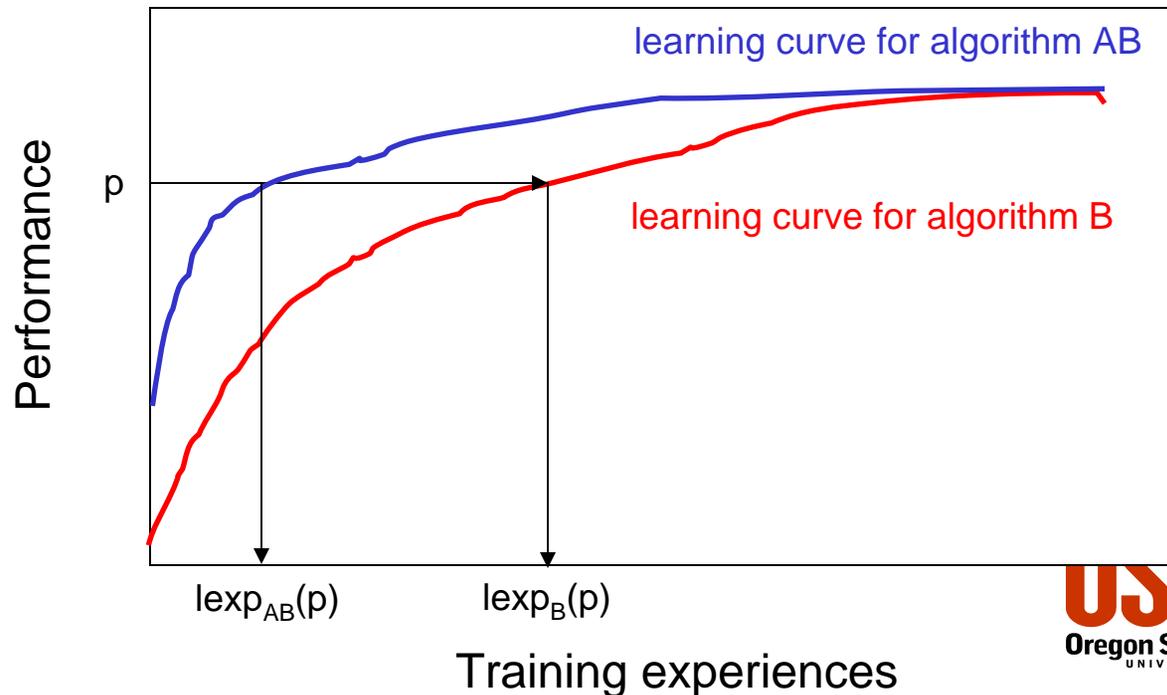
Learning Speedup and Relative Reduction

- Let
 - algorithm B = learning without transfer
 - algorithm AB = learning with transfer
- The speedup in the number of training examples required to reach target performance level P is

$$\text{speedup}(p) = \frac{\text{lexp}_B(p)}{\text{lexp}_{AB}(p)}$$

- The relative reduction in amount of training required is

$$\text{RR}(p) = \frac{\text{lexp}_B(p) - \text{lexp}_{AB}(p)}{\text{lexp}_B(p)} = 1 - \frac{\text{lexp}_{AB}(p)}{\text{lexp}_B(p)}$$



- The asymptotic value of the ratio

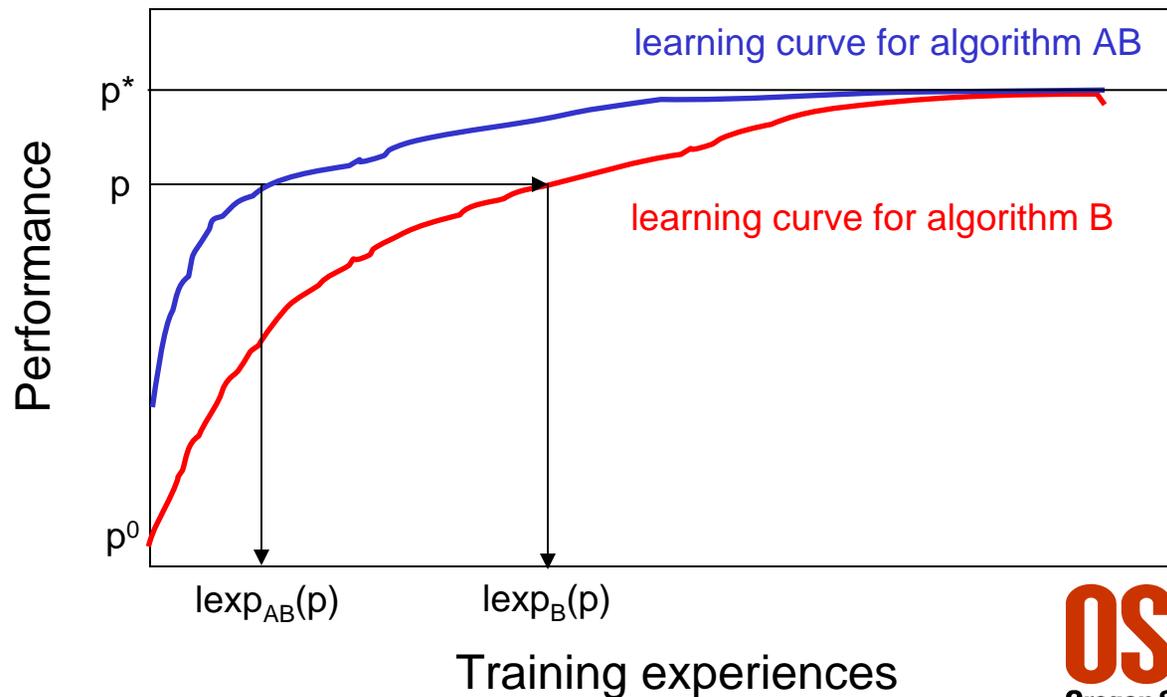
$$\lim_{p \rightarrow p^*} \frac{\text{lexp}_B(p)}{\text{lexp}_{AB}(p)}$$

is known as the statistical efficiency of algorithm AB relative to algorithm A. It is a standard measure employed in statistics to compare *consistent* learning algorithms – that is, algorithms that reach optimal performance asymptotically

Integrated Speedup

- Goal: obtain a single number that summarizes the relative performance.
- Solution: Integrate this speedup p^0 up to the asymptote for the task:

$$\frac{1}{p^* - p^0} \int_{p^0}^{p^*} \text{speedup}(p) dp = \frac{1}{p^* - p^0} \int_{p^0}^{p^*} \frac{\text{lexp}_B(p)}{\text{lexp}_{AB}(p)} dp$$

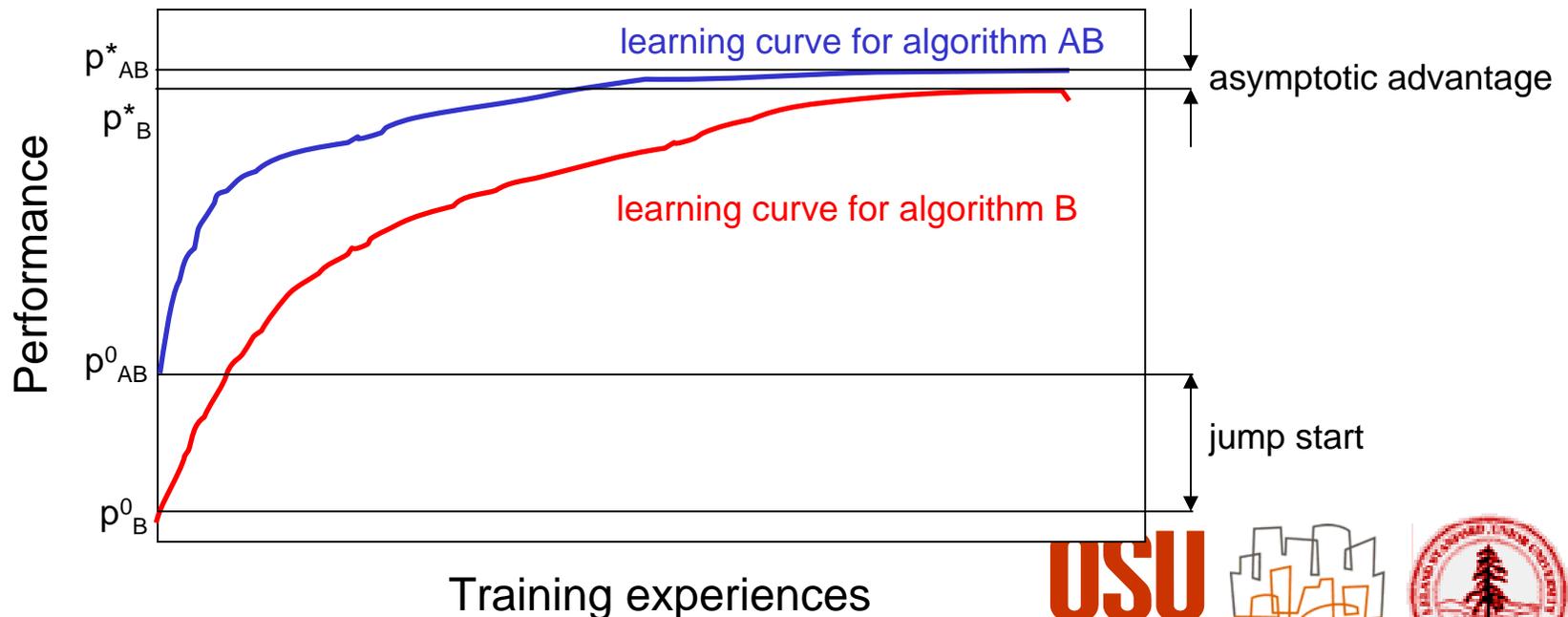


- The ratio will be infinite if the AB curve has a Y-axis intercept above the B curve intercept
 - at such points, we have infinite speedup
- The ratio will be infinite if the AB curve has an asymptote above the B curve asymptote
 - for levels $p >$ the B curve asymptote, we have infinite speedup

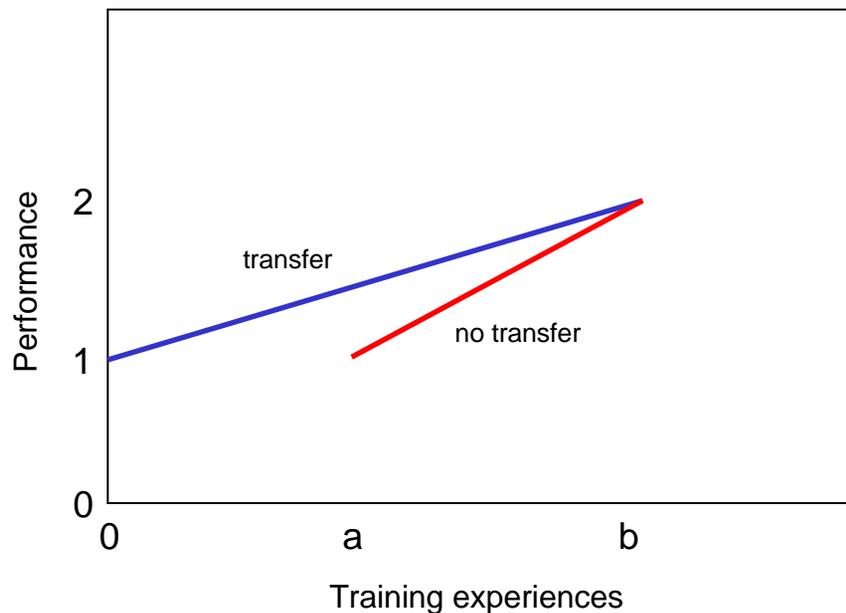
Partial Solution: Remove the Problem Regions

- Define three quantities:
 - “jump start”: $p_{AB}^0 - p_B^0$
 - “asymptotic advantage”: $p_{AB}^* - p_B^*$
 - “integrated speedup”:

$$\frac{1}{p_B^* - p_{AB}^0} \int_{p_{AB}^0}^{p_B^*} \text{speedup}(p) dp = \frac{1}{p_B^* - p_{AB}^0} \int_{p_{AB}^0}^{p_B^*} \frac{\text{lexp}_B(p)}{\text{lexp}_{AB}(p)} dp$$



- The integrated speedup can still be infinite
- At p_{AB}^0 , the speedup is usually infinite
 - although this occurs only at a point, it can cause the integral to diverge
 - example:



$$\int_1^2 \frac{(p-1)(b-a) + a}{(p-1)b} dp =$$

$$\int_1^2 \left[\frac{b-a}{b} + \frac{a}{b} \frac{1}{p-1} \right] dp =$$

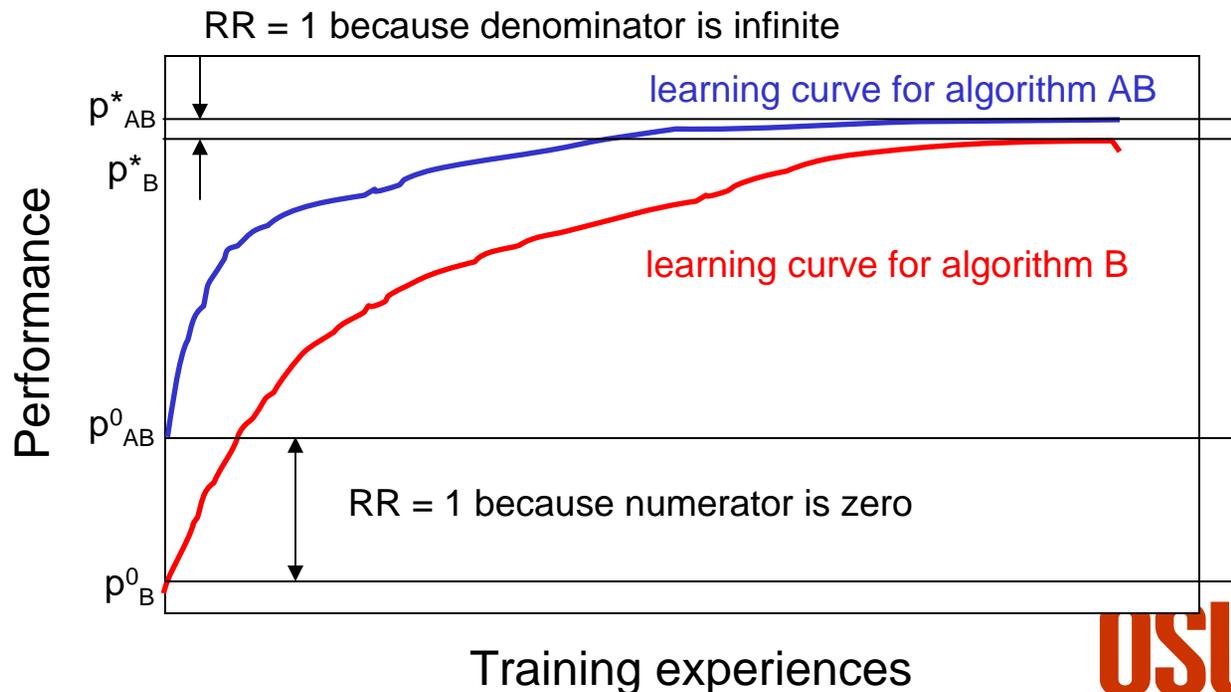
$$\frac{b-a}{b} + \frac{a}{b} \ln p \Big|_0^1 = \infty$$

Solution: Use Relative Reduction Instead

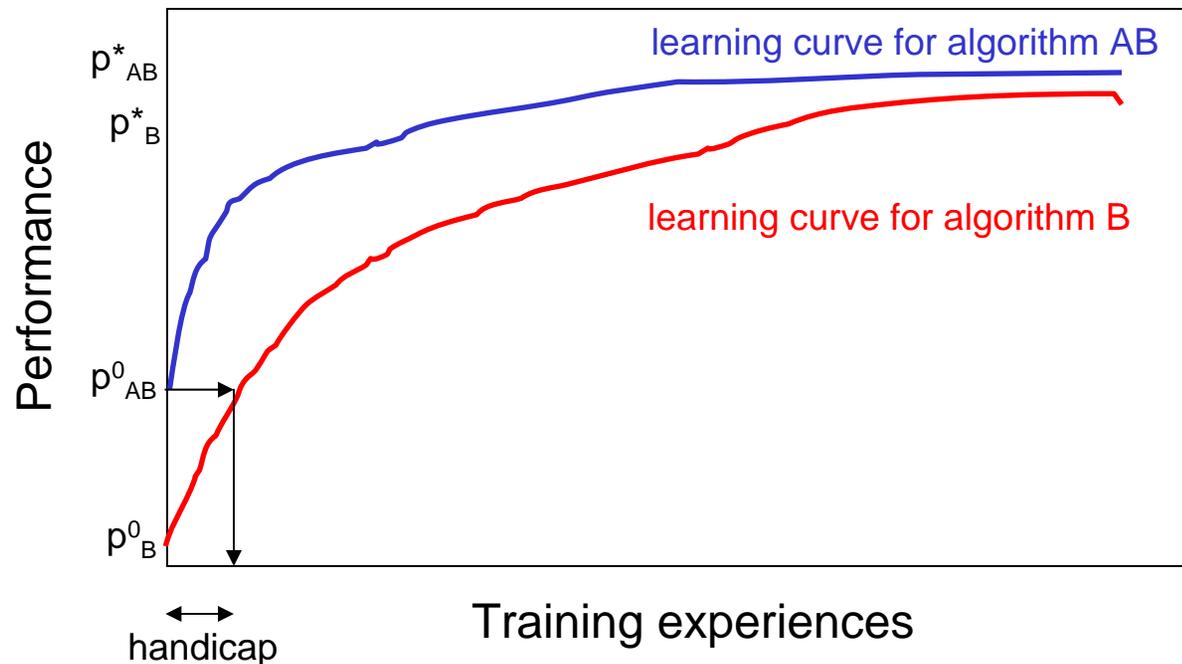
- “average relative reduction”:

$$\frac{1}{p_{AB}^* - p_B^0} \int_{p_B^0}^{p_{AB}^*} RR(p) dp = \frac{1}{p_{AB}^* - p_B^0} \int_{p_B^0}^{p_{AB}^*} 1 - \frac{\text{lexp}_{AB}(p)}{\text{lexp}_B(p)} dp$$

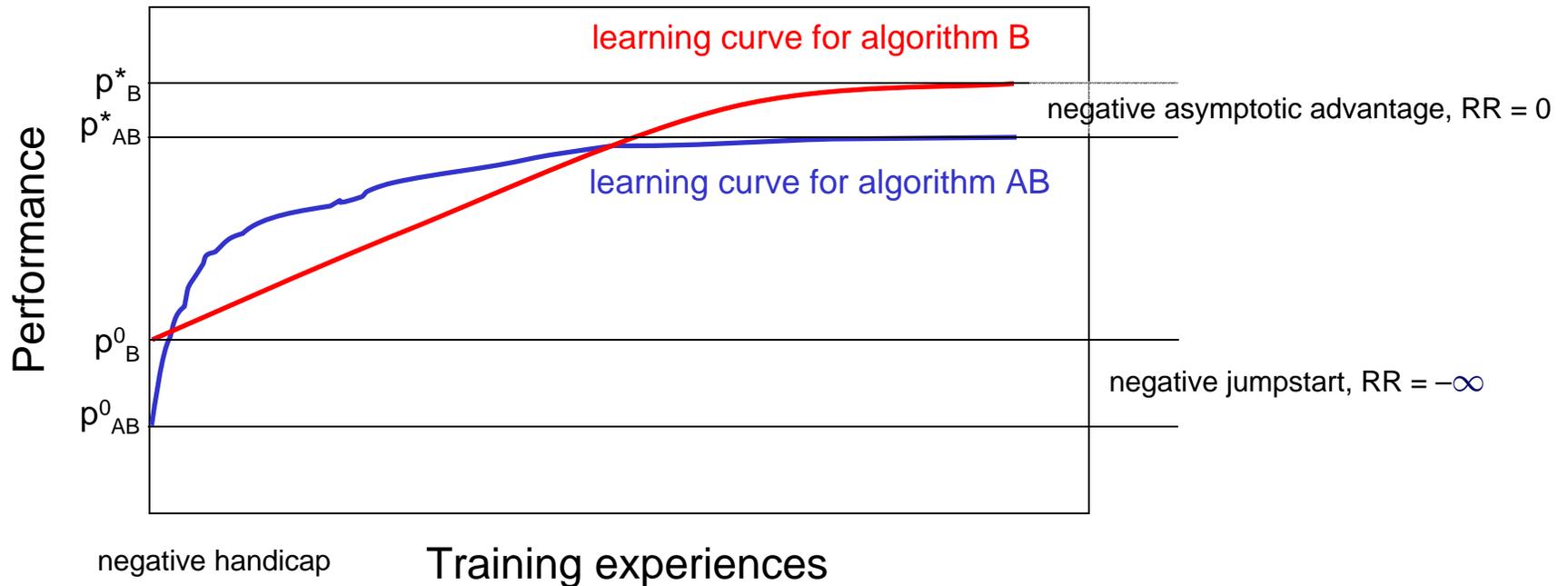
- avoids (almost) all of the division-by-zero problems
- define integral to run from $\min(p_B^0, p_{AB}^0)$ to $\max(p_B^*, p_{AB}^*)$



- “handicap”: How much training does it take B to reach performance level p^0_{AB} ? (i.e., to overcome the jumpstart)

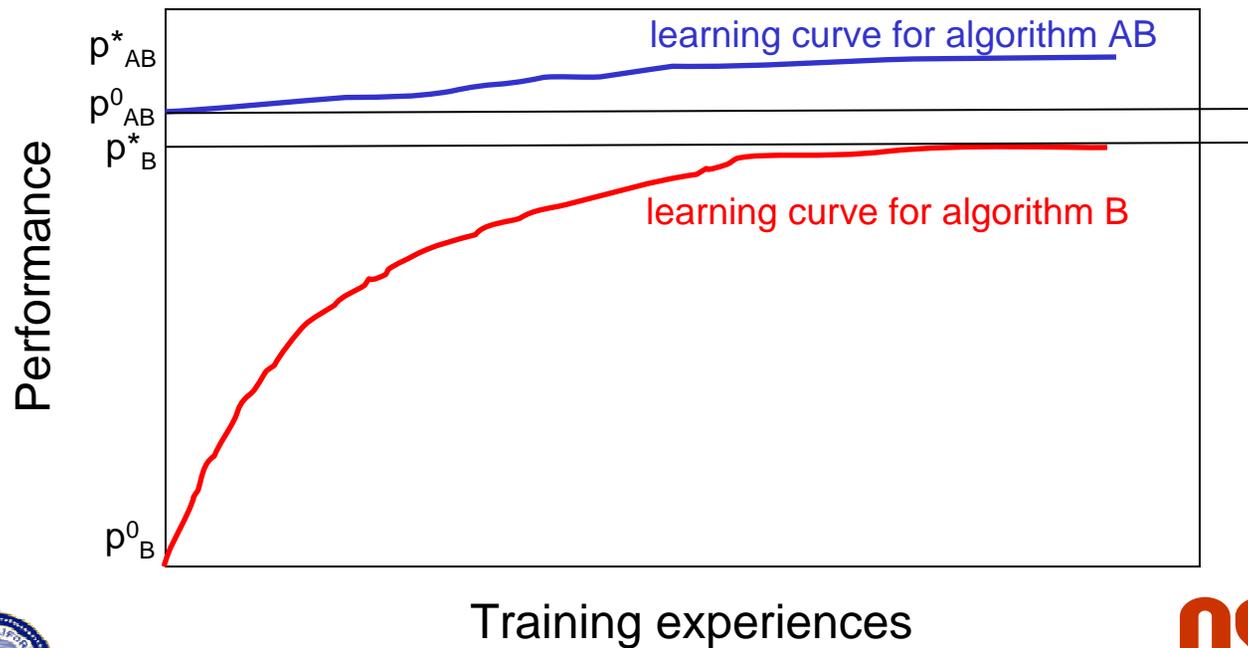


- Crossing curves: all 3 auxiliary measures can be negative



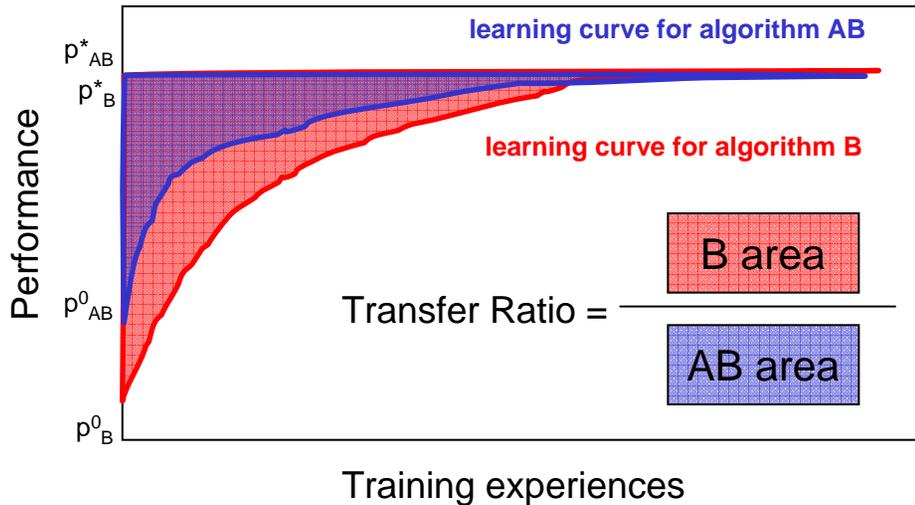
- Average RR will be $-\infty$. This will be rare, but it gives the no-transfer case the benefit.

- Transfer curve can be completely above the no-transfer curve:
 - large jumpstart
 - modest asymptotic advantage
 - infinite handicap
 - $RR = 1$ everywhere

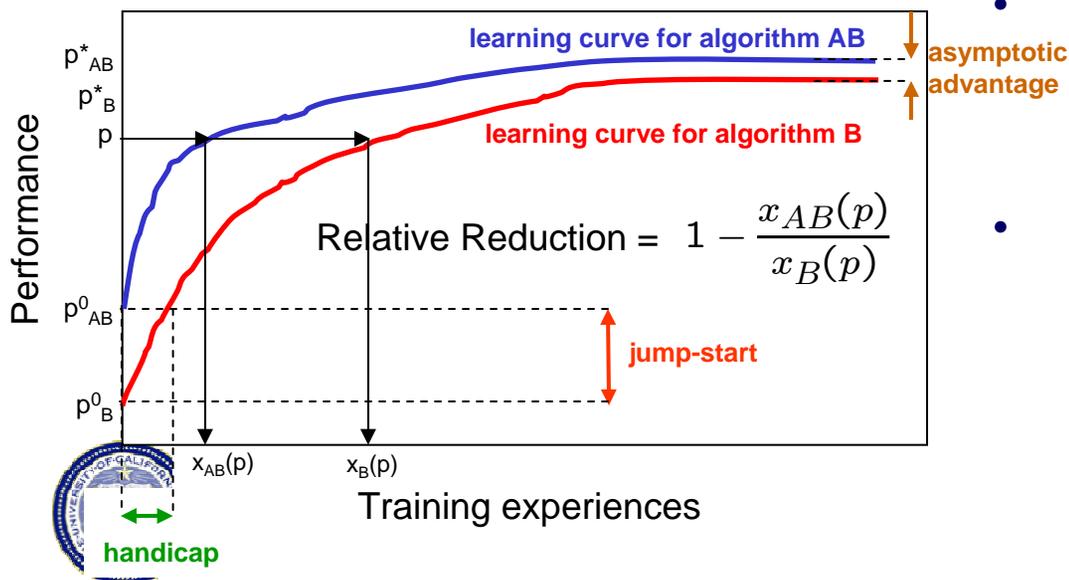


- The algorithm(s) can attain performance level p multiple times. $\text{lexp}(p)$ must be defined as the minimum of such occurrences
- The RR integral can be < 0 , which indicates that transfer learning was worse than no-transfer
- In practice, each learning curve must be approximated, and because it is measured at fixed training experiences N_1, N_2, \dots , it is unlikely that the measured performance will exactly match a given target level p . Hence we propose approximating the learning curve by a piece-wise linear curve and interpolating p and N linearly to compute $\text{lexp}(p)$. With this approximation, the integral can be computed exactly (over a series of trapezoids defined by $\text{perf}(N_1), \text{perf}(N_2), \dots$).
- Of course the measured learning curve is only an estimate of the true curve. We could generate bootstrap learning curves (via bootstrap replicates of the test data), compute the metrics for each curve, and then compute and report confidence intervals for the metrics.

- Four proposed metrics:
 - average relative reduction in training time (sample size, number of training experiences)
 - jumpstart (initial advantage of transfer algorithm)
 - handicap (how long it takes the no-transfer algorithm to overcome the jumpstart)
 - asymptotic advantage (how much better the transfer learning algorithm does in the limit of large sample sizes)



- Problems:
 - Does not give instantaneous measure
 - Assumes $p_{AB}^* = p_B^*$ (gives infinite ratio otherwise)
 - Hard to explain (unclear units)
- Solution:
 - Compute instantaneous measure, then integrate
 - Integrate along the performance axis
 - Measure reduction in amount of training needed ("amount reduced by 75%")



- Average relative reduction:

$$\frac{1}{p_{AB}^* - p_B^0} \int_{p_B^0}^{p_{AB}^*} 1 - \frac{x_{AB}(p)}{x_B(p)} dp$$

- Auxiliary metrics:

- jump start: $p_{AB}^0 - p_B^0$
- asymptotic advantage: $p_{AB}^* - p_B^*$
- handicap: $x_B(p_{AB}^0)$