Click here Slide(pdf)) a way to find "similiar" sets. LSH: Locality-Sensitive Hashing upside: only a small faction of points are ever examined downside: exists false negatives 1. hash shingles to a few bytes fisrt. Shingling 2. Compute Sim(C1, C2) Find a hash function h(x) make that: second. Min-hashing The function to our taste is Min-hashing Steps for similiar docs Then we get signatures. Focus on pairs of signatures likely to be from similiar documents There uses bands and rows to reduce the error. third. LSH Use hashing to find candidate pairs of similarity >= s