## The Curse of Dimensionality

After **overfitting**, the biggest problem in machine learning is the curse of dimensionality. It refers to the fact that many algorithms that work fine in low dimensions become intractable when the input is high-dimensional. Generalizing correctly becomes exponentially harder as the dimensionality (number of features) of the examples grows, because a fixed-size training set covers a swindling fraction of the input space. Even with a moderate dimension of 100 and a huge training n set of a trilling examples, the latter covers only a fraction of about $10^{-18}$ of the input space. This is what makes machine learning both necessary and hard.

Our intuitions, which come from a three-dimensional world, often do not apply in high-dimensional ones. In high dimensions, most of the mass of a multivariate Gaussian distribution is not near the mean, but in increasingly distant "shall" around it. Building a classifier in two or three dimensions is easy; we can find a reasonable frontier between examples of different classes just by visual inspection. But in high dimensions it's hard to understand what is happening. This in turn makes it difficult to design a good classifier. Naively, one might think that gathering more features never hurts, since at worst they provide no new information about the class. But in fact their benefits may be outweighed by the curse of dimensionality.