Assignment 2:

1) Possible to do slightly better than decision trees for tumor, but not by much

2) Overfitting with small K, _maybe_ (depends on noise, intrinsic dimensionality of data)

3) KNN is quite slow if distance evaluation is slow

4) In large dimensions, vectors start to look <u>very</u> <u>orthogonal to one another.</u>

5) Lower dimensionality makes it harder to overfit

6) Maybe, if you are concerned about overfitting or noise

How do the answers to 2) and 4) connect to one another? What about 2) and 6), and 4) and 5)

Agaricus-Lepiota has higher dimensionality than primary-tumor, but higher accuracy as well ?!

# Quiz

1) Why is leave-one-out cross validation practical for K-NN classifiers and impractical for other settings?

2) For any bootstrap sample, what's the expected number of times any element of the original sample will be chosen?

In the limit of an infinitely large sample, what is the probability that any given element of the sample is selected?

Probability that a sample isn't selected in one of the $n$ draws: $\dfrac{n-1}{n}$

Probability it isn't selected in any:

$$\lim_{n \to \infty} \left(\frac{n-1}{n}\right)^n = \left(1 - \frac{1}{n}\right)^n = e^{-1}$$

Probability it is selected at least once: $1 - \dfrac{1}{e} \approx 0.63$

# Precision, Recall, Sensitivity, Oh my!

Confusion matrix

$$Acc = \frac{A+D}{A+B+C+D}$$

| | $\hat{y}=0$ | $\hat{y}=1$ |
|---|---|---|
| $y=0$ | A | B |
| $y=1$ | C | D |

$$T = C + D$$
$$S = B + D$$
$$I = D$$

$$\text{Precision} = \frac{I}{S} = \frac{D}{B+D}$$

$$\neg I = A$$
$$\neg S = A + C$$

$$\text{Recall} = \frac{I}{T} = \frac{D}{C+D}$$

$$\text{Sensitivity} = \text{Recall} = \frac{D}{C+D}$$

$$\text{Specificity} = \frac{\neg I}{\neg S} = \frac{A}{A+C}$$

# Cross Validation

- Validation data is used inefficiently

  - Let's partition training data into K subsets, then create K training-validation splits

  - The average of the validation error is a better estimate of test error

  - "K-fold cross validation"

  - Use this to decide on hyperparameters

- What if K increases to $\infty$?

- leave-one-out CV
- Only useful when not horribly inefficient

- KNN

**Algorithm 9** KNN-Train-LOO(D)

1: $err_k \leftarrow 0, \forall 1 \leq k \leq N-1$      // $err_k$ stores how well you do with $k$NN
2: **for** $n = 1$ **to** $N$ **do**
3:     $S_m \leftarrow \langle ||x_n - x_m||, m\rangle, \forall m \neq n$      // compute distances to other points
4:     $S \leftarrow \text{SORT}(S)$      // put lowest-distance objects first
5:     $\hat{y} \leftarrow 0$      // current label prediction
6:     **for** $k = 1$ **to** $N-1$ **do**
7:       $\langle dist, m\rangle \leftarrow S_k$
8:       $\hat{y} \leftarrow \hat{y} + y_m$      // let $k$th closest point vote
9:       **if** $\hat{y} \neq y_m$ **then**
10:         $err_k \leftarrow err_k + 1$      // one more error for $k$NN
11:       **end if**
12:     **end for**
13: **end for**
14: **return** $\text{argmin}_k err_k$      // return the $K$ that achieved lowest error

# Bootstrap Resampling (Statistical Testing)

- When is 8% really greater than 7.9%?
  - fundamental new idea: error estimation matters

- Null hypothesis testing asks:

  - How likely is it that we saw this difference because of chance?

  - If we could repeat the experiment infinitely many times, how often would we have made a mistake?

- Instead of reporting "8% on our 'population'", we report (an estimate of) the distribution by rerunning our experiment on (many estimates of, resamples of) the population

---

**Algorithm 10** BOOTSTRAPEVALUATE($y$, $\hat{y}$, $NumFolds$)

```
1:  scores ← [ ]
2:  for k = 1 to NumFolds do
3:      truth ← [ ]                                    // list of values we want to predict
4:      pred ← [ ]                                     // list of values we actually predicted
5:      for n = 1 to N do
6:          m ← uniform random value from 1 to N       // sample a test point
7:          truth ← truth ⊕ yₘ                          // add on the Truth
8:          pred ← pred ⊕ ŷₘ                            // add on our prediction
9:      end for
10:     scores ← scores ⊕ F-SCORE(truth, pred)          // evaluate
11: end for
12: return (MEAN(scores), STDDEV(scores))
```