# The k-NN classifier
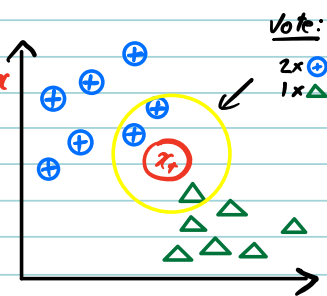
<u>Assumption</u>: Similar points share similar labels

<u>Classification Rule</u>: For a test input $x_t$ assign the <span style="color:red">most common</span> label among its $k$ most similar <span style="color:red">training</span> inputs.

<u>Formally</u>: $D = \{(x_1, y_1), \ldots, (x_n, y_n)\}$ training data. Test point $x$

Let $S_x \subseteq D$ such that $|S_x| = k$
and $\forall (x', y) \in D \backslash S_x$   $dist(x', x) \geq \max\limits_{(x'', y'') \in D \backslash S_x} dist(x'', x)$

$h(x) = \text{mode}\{y': (x', y') \in S_x\}$

Vote:
$2 \times \oplus$
$1 \times \triangle$

<u>Protip</u>: In case of a draw decide by reducing $k$ by 1, until you reach a unique mode.

<u>Training Error</u>: Leave-One-Out (LOO) estimate: Take each training point out and estimate its label, pretending it was a test point. (i.e. a point cannot be its own neighbor)

<u>What distance function should we use?</u>

Common choice: Minkowki's distance:   $$dist(x, z) = \left(\sum_\alpha |[x]_\alpha - [z]_\alpha|^p\right)^{1/p} \quad \text{for} \quad p > 0$$

special case:  $p = 2 \leftarrow$ Euclidean distance
$p = 1 \leftarrow$ Manhattan distance

<u>Quiz</u>: What if $p \to \infty$ or $p \to 0^+$ ?
How does $k$ affect the outcome? How does the classifier behave as $k=1$, or $k=n$?

## Bayes Optimal Classifier

Your data $D$ is drawn from some distribution $(x, y) \sim P(x, y)$. Also: $P(x, y) = P(y|x)P(x)$
Assume you knew $P(y|x)$ (you never do, but just for the sake of the argument).
For some test $x$ what label would you predict?

<span style="color:blue">The most likely label:</span> $h_{opt}(x) = \text{argmax}_y P(y|x)$

What is the expected error of the BOC? Let $y^* = h_{opt}(x)$   $\varepsilon = 1 - P(y^*|x)$

<span style="color:orange">The probability that $x$ does not have the most likely label.</span>

You can never do better than the BOC!

# Asymptotic error bound for 1-NN
## (Cover and Hart 1967)

Quiz: 1. You have a coin that shows head with probability, $p$.
If you throw it twice, what is the probability $q$ that both throws lead to *different* outcomes?

2. Show that $q \leq 2(1-p)$.
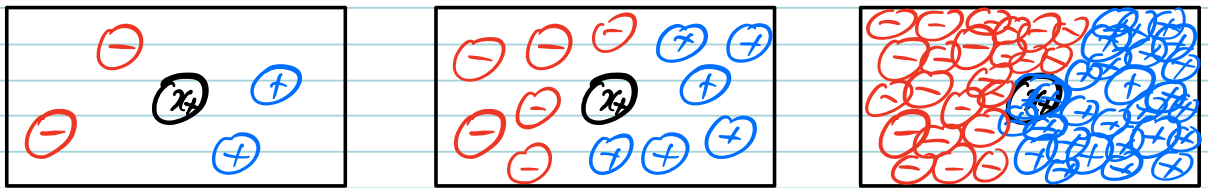
Back to 1-NN. We want to prove that the expected 1-NN test error is less than $2\times$ the BOC error, as $n \to \infty$. (For binary classification.)

Argument: Let $x$ be the test point and $\hat{x}$ be its nearest neighbor.

Claim 1: As $n \to \infty$, $\text{dist}(x, \hat{x}) \to 0$   ← i.e. The nearest neighbor becomes infinitely close.

Claim 2: As $\text{dist}(x, \hat{x}) \to 0$, $\hat{x} \to x$   ← i.e. In fact, the nearest neighbor becomes identical to $x$.

(See Covert & Hart for proof.)



Assume for $x_t$ the label $y^*$ is most likely. Let $\rho = P(y^* | x_t)$

The BOC would predict $y^*$, and be wrong with probability $\varepsilon_{BOC} = 1 - \rho$.

What is the error of 1-NN as $n \to \infty$?

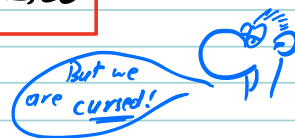1-NN is wrong if the labels of $x$ and $\hat{x}$ are *different*.
By claim 2 we have $\hat{x} \to x$. And $\rho(y^* | \hat{x}) = \rho(y^* | x) = \rho$
Both points $x$ and $\hat{x}$ could take on label $y^*$ with prob. $\rho$, and not with $(1-\rho)$.

Remember Quiz 2. Regard both points as the same coin tossed twice.
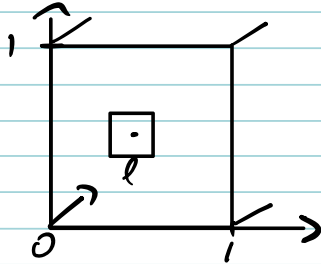They disagree with probability $2\rho(1-\rho) \leq 2(1-\rho) = 2\varepsilon_{BOC}$

$\Rightarrow \boxed{\varepsilon_{1-NN} \leq 2\varepsilon_{BOC}}$   as $n \to \infty$

Yay!

But we are cursed!

# Curse of Dimensionality

Assume $x_i \in [0,1]^d$ (i.e. the $d$ dimensional unit hypercube).
All data is drawn uniformly at random.
Let $k = 10$.

Let $l$ be the edge length of the smallest hypercube that contains all $k$ nearest neighbors of a test point $x$.

$$l^d \approx \frac{k}{n} \implies l \approx \left(\frac{k}{n}\right)^{1/d}$$

↑ volume of mini cube containing the $k$ neighbors

↑ Total volume of hypercube $[0,1]^d$ is $1^d = 1$
$\frac{k}{n}$ is the fraction that $k$ points take up. (because points are uniformly sampled)

If $n = 1000$ how big is $l$?

| $d$ | 2 | 10 | 100 | 1000 |
|---|---|---|---|---|
| $l$ | 0.1 | 0.63 | 0.955 | 0.9954 |

Almost the entire space is needed to fit 10 nearest neighbors.

This means nearest neighbors are not similar, violating the k-NN assumption!

How many points would we need for $l$ to be small?
Fix $l = 0.1$

$$l^d = \frac{k}{n} \implies n = k\left(\frac{1}{l}\right)^d = k \, 10^d$$ ← grows exponentially with $d$!

## Rescue to the curse:

Data can have structure:

- Data can lie on intrinsically low dimensional subspaces or sub-manifolds.

- Data can be clustered (very non-uniform).