# Machine Learning from Data

Lecture 16: Spring 2021

# Today's Lecture

- Similarity
- Nearest Neighbor

FOUNDATIONS

Theory.
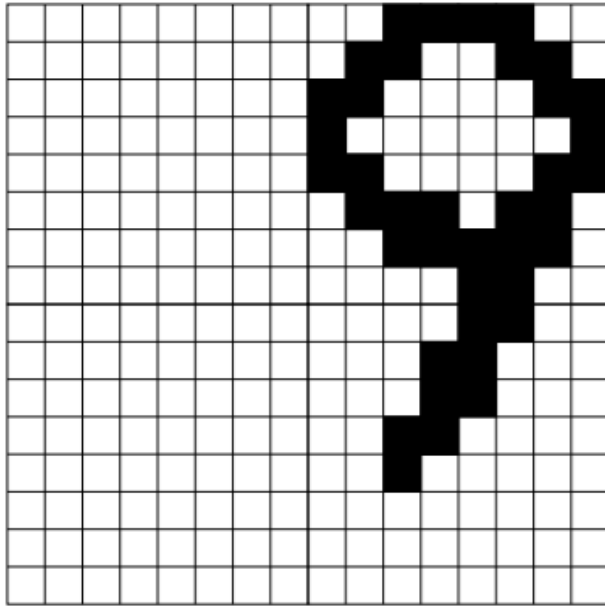
MANOHORSE

*frame work*
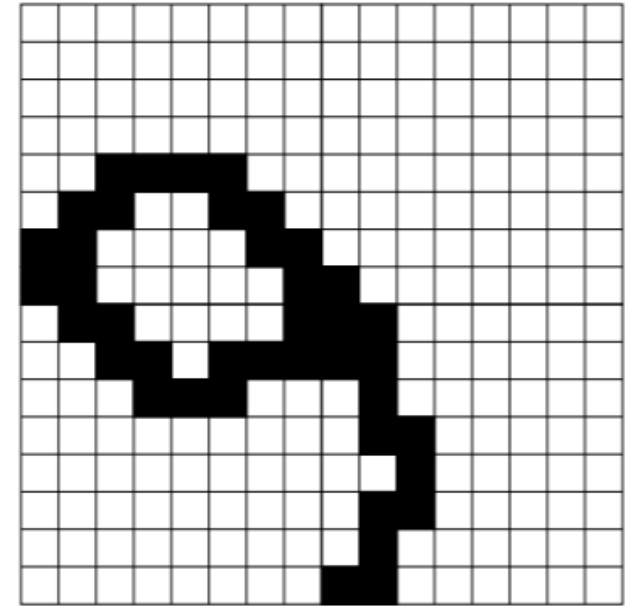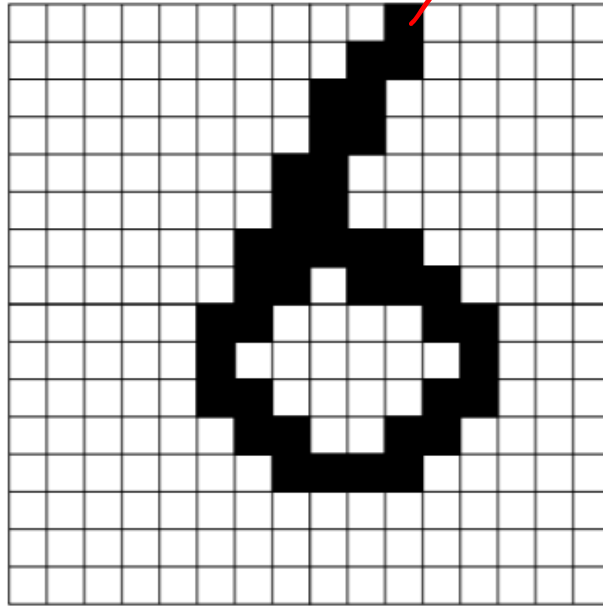
Ask a 5-year-old what is this?

NN rule

9

6 → counting* 9

Capture essense of features w.r.t. the problem
Created/Constructed

$$x \quad x'$$

$$d(x, x') = \|x - x'\|$$

Euclidean

Similarity measure

Test point $\rightarrow x$

NN - Algorithm

$x_1^0$

$x_{[2]} \quad x_{[4]}$

$x_2$ $\quad x_{[1]} \quad x_{[3]}$

$$g(x) = \text{sign}[y_{[1]}]$$

$x_3^0$

$x_4^0$
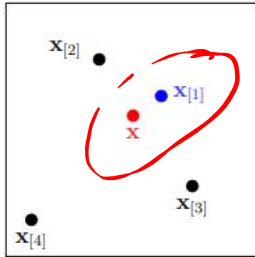
Voronoi

Test 'x' is classified using its nearest neighbor.
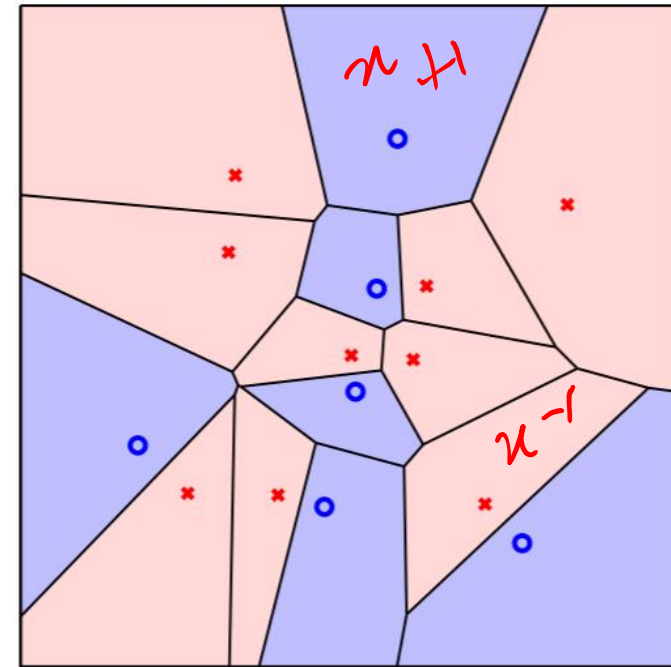
$$d(\mathbf{x}, \mathbf{x}_{[1]}) \leq d(\mathbf{x}, \mathbf{x}_{[2]}) \leq \cdots \leq d(\mathbf{x}, \mathbf{x}_{[N]})$$

$x_{[2]}$
$x_{[1]}$
x
$x_{[3]}$
$x_{[4]}$

$$g(\mathbf{x}) = y_{[1]}(\mathbf{x})$$

**No training needed!**

$E_{\text{in}} = 0$

$E_{in}$ VS $E_{out}$

$\varkappa H$

$\varkappa -1$

Nearest neighbor Voronoi tesselation

$$d_{VC} = \infty \qquad E_{in} = 0$$

$E_{in}$ and $E_{out}$

$\rightarrow$ 5-year old.

$$\underline{\text{Theorem}} \quad : \quad \boxed{E_{out} \leq 2 E_{out}^{*}}$$

$\begin{cases} \text{w. h. p} \text{ --} \\ \text{\& sufficiently} \\ \text{large } N \text{ --} \end{cases}$

$E_{out}^{*} \longrightarrow$ $\underline{\text{Optimal } E_{out}} \longrightarrow$ The best possible out-of-sample error. for a given problem.

How small can $E_{out}$ get ?

**Proof :** $f(x) \Rightarrow \underline{\underline{\pi(x)}} = P(y=+1|x)$

$\underline{\text{Optimal classifier :}}$

$\left. \begin{array}{ll} +1 & \text{if} \quad \pi(x) \geqslant \frac{1}{2} \\ -1 & \text{if} \quad \pi(x) < \frac{1}{2} \end{array} \right\}$

$E_{out}^{*}(x)$
$= \int dx \, P(x) \overline{E}_{out}^{*}(x)$

$\underline{\underline{E_{out}^{*}(x)}} = \left\{ \begin{array}{ll} 1 - \pi(x), & \pi(x) \geqslant \frac{1}{2} \\ \\ \pi(x), & \pi(x) < \frac{1}{2} \end{array} \right\}$

$\underline{\underline{E_{out}^{*}(x)}} = \min\left(\pi(x), 1 - \pi(x)\right) = \eta(x)$
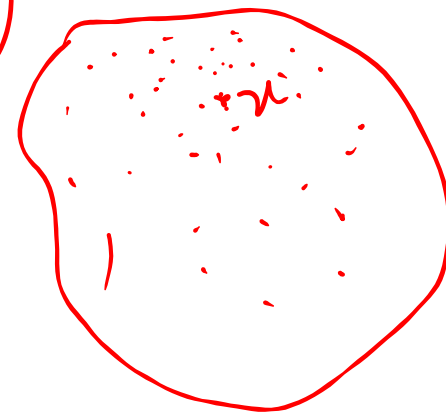
1) $\pi(x)$ is continuos

2) Data generation

$x_{[i]}$     $N \longrightarrow \infty$   $\Rightarrow$   $\underline{x_{[i]}} \longrightarrow x$

$$\text{ProbError}[x] = P\left[y_{[i]} \neq y\right]$$

$y_{[i]} <$

$$= P\left[y_{[i]} = +1, y = -1\right] + P\left[y_{[i]} = -1, y = +1\right]$$

$$= \pi(x_{[i]})(1 - \pi(x)) + (1 - \pi(x_{[i]}))\pi(x)$$

N is sufficiently large:

$$x_{[1]} \longrightarrow x$$

$$\pi(x_{[1]}) \longrightarrow \pi(x)$$

$$\text{ProbError}[x] = \pi(x)(1 - \pi(x)) + (1 - \pi(x))\pi(x)$$

$$= 2\,\pi(x)(1 - \pi(x))$$

$$\leq 2\,\eta(x) = 2E_{out}^{*}(x)$$

$$E_{out}(x) \leq 2E_{out}^{*}(x)$$

# Proving $E_{\text{out}} \leq 2E_{\text{out}}^*$

$$\pi(\mathbf{x}) = \mathbb{P}[y = +1 | \mathbf{x}].$$

Assume $\pi(\mathbf{x})$ is continuous and $\mathbf{x}_{[1]} \xrightarrow{N \to \infty} \mathbf{x}$. Then $\pi(\mathbf{x}_{[1]}) \xrightarrow{N \to \infty} \pi(\mathbf{x})$.

$$
\begin{aligned}
\mathbb{P}[g_N(\mathbf{x}) \neq y] &= \mathbb{P}[y = +1, y_{[1]} = -1] + \mathbb{P}[y = -1, y_{[1]} = +1], \\
&= \pi(\mathbf{x}) \cdot (1 - \pi(\mathbf{x}_{[1]})) + (1 - \pi(\mathbf{x})) \cdot \pi(\mathbf{x}_{[1]}), \\
&\to \pi(\mathbf{x}) \cdot (1 - \pi(\mathbf{x})) + (1 - \pi(\mathbf{x})) \cdot \pi(\mathbf{x}), \\
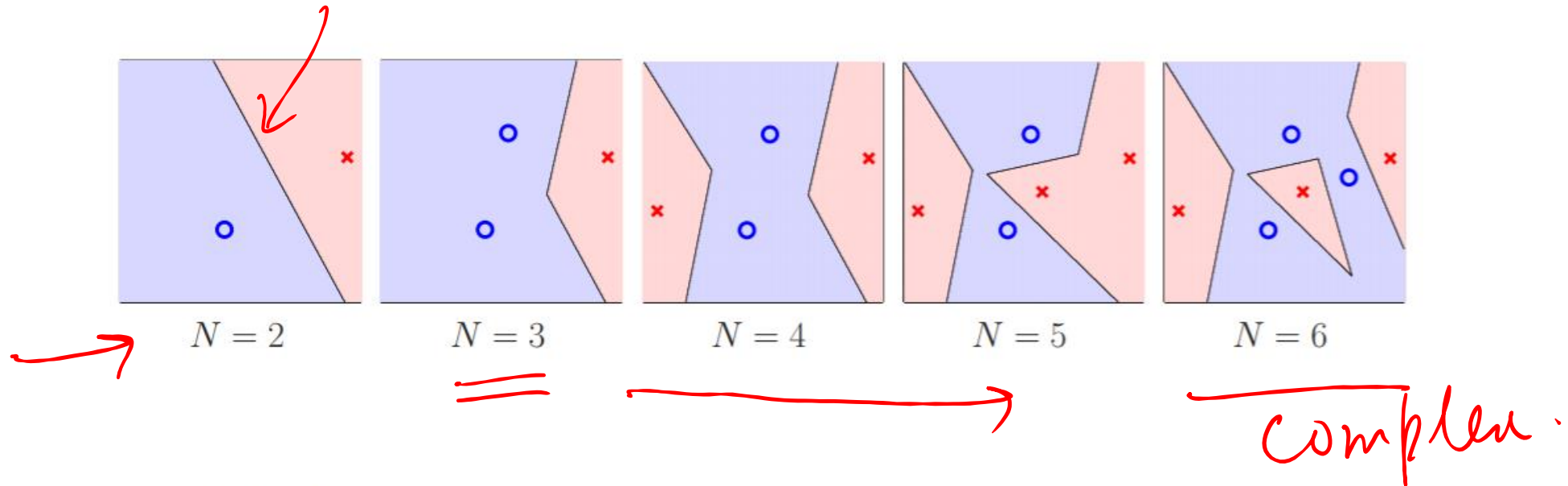&= 2\pi(\mathbf{x}) \cdot (1 - \pi(\mathbf{x})), \\
&\leq 2\min\{\pi(\mathbf{x}), 1 - \pi(\mathbf{x})\}.
\end{aligned}
$$

The best you can do is

$$E_{\text{out}}^*(\mathbf{x}) = \min\{\pi(\mathbf{x}), 1 - \pi(\mathbf{x})\}.$$

# Nearest Neighbor 'Self-Regularizes'



$N = 2$      $N = 3$      $N = 4$      $N = 5$      $N = 6$

*complex.*

A simple boundary is used with few data points.

A more complicated boundary is possible *only* when you have more data points.

regularization guides you to simpler hypotheses when data quality/quantity is lower.
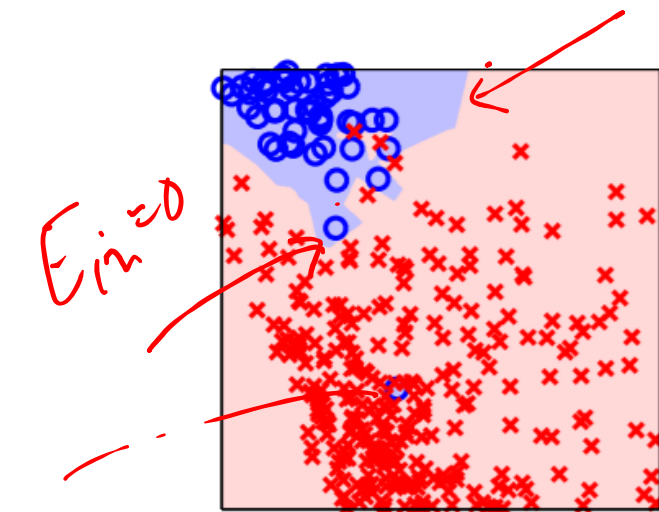
# $k$-Nearest Neighbor

$$g(\mathbf{x}) = \text{sign}\left(\sum_{i=1}^{k} y_{[i]}(\mathbf{x})\right).$$
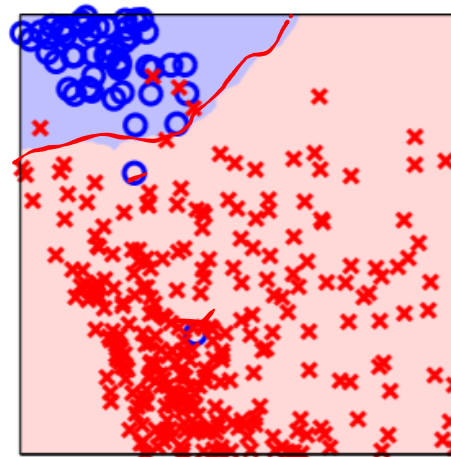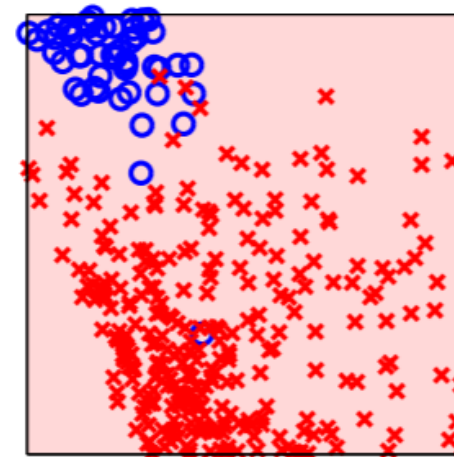
($k$ is odd and $y_n = \pm 1$).

Majority

k-regularization



$E_{in} = 0$

1-NN rule ✓          21-NN rule          127-NN rule ✓

overfit                                    Underfit

# The Role of $k$

$k$ determines the tradeoff between fitting the data and overfitting the data.

**Theorem.** For $N \to \infty$, if $k(N) \to \infty$ and $k(N)/N \to 0$ then,

$$E_{\text{in}}(g) \to E_{\text{out}}(g) \quad \text{and} \quad E_{\text{out}}(g) \to E_{\text{out}}^*.$$

For example $k = \left\lceil \sqrt{N} \right\rceil$.

*(handwritten annotations)* choose $k$ → CV    $k/N$    $\lambda(n)$    converges    $K = 3$
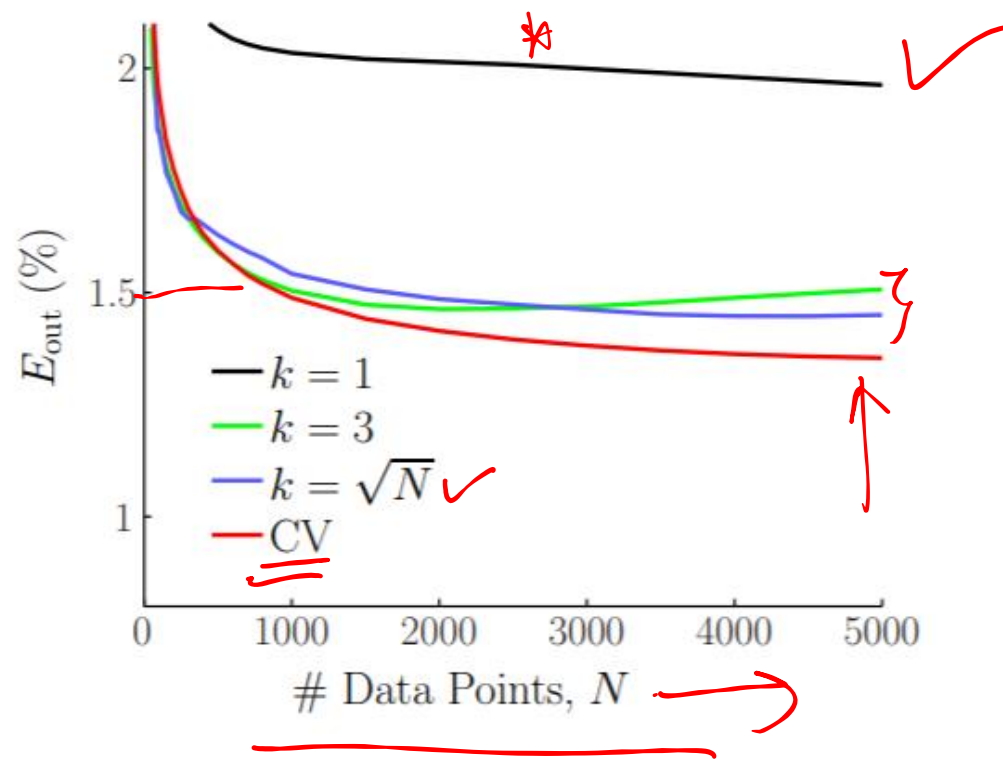
# 3 Ways To Choose $k$

1. $k = 3$.

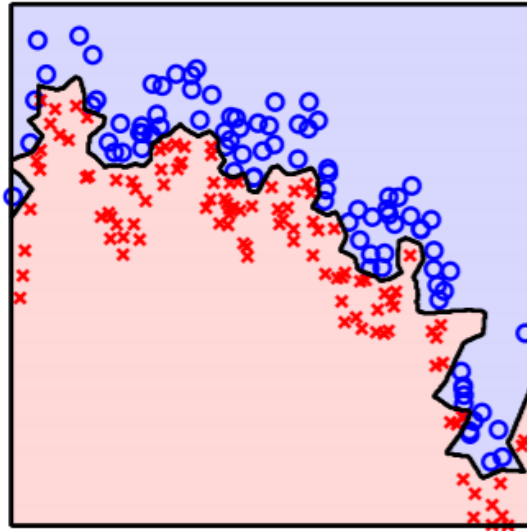2. $k = \lceil \sqrt{N} \rceil$.

3. Validation or cross validation:

   $k$-NN rule hypotheses $g_{\bar{k}}$ constructed on training
   set, tested on validation set, and best $k$ is picked.

# Nearest Neighbor is Nonparametric

### NN-rule



### Linear Model (Parametric)



→ line

no parameters

expressive/flexible

$g(\mathbf{x})$ needs data

generic, can model anything

$(d+1)$ parameters

rigid, always linear

$g(\mathbf{x})$ needs only weights ← keep

specialized

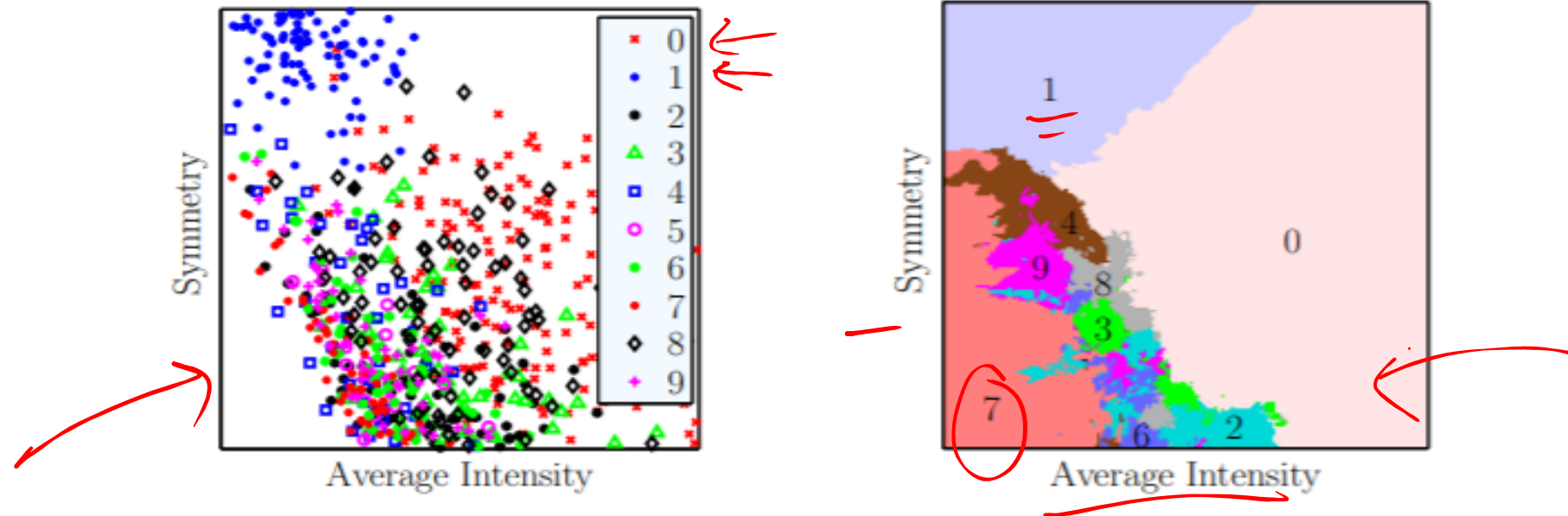# Nearest Neighbor Easily Extends to Multiclass



**Confusion Matrix**

| True | Predicted | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | |
| 0 | **13.5** | 0.5 | 0.5 | 1 | 0 | 0.5 | 0 | 0 | 0.5 | 0 | 16.5 |
| 1 | 0.5 | **13.5** | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 14 |
| 2 | 0.5 | 0 | **3.5** | 1 | 1 | 1.5 | 1 | 1 | 0 | 0.5 | 10 |
| 3 | 2.5 | 0 | 1.5 | **2** | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 1 | 9.5 |
| 4 | 0.5 | 0 | 1 | 0.5 | **1.5** | 0.5 | 1 | 2 | 0 | 1.5 | 8.5 |
| 5 | 0.5 | 0 | 2.5 | 1 | 0.5 | **1.5** | 1 | 1 | 0 | 0.5 | 7.5 |
| 6 | 0.5 | 0 | 2 | 1 | 1 | 1 | **1** | 1 | 0 | 1 | 8.5 |
| 7 | 0 | 0 | 1.5 | 0.5 | 1.5 | 0.5 | 1 | **3** | 0 | 1 | 9 |
| 8 | 3.5 | 0 | 0.5 | 1 | 0.5 | 0.5 | 0.5 | 0 | **0.5** | 1 | 8 |
| 9 | 0.5 | 0 | 1 | 1 | 1 | 0.5 | 1 | 1 | 0.5 | **2** | 8.5 |
| | 22.5 | 14 | 14 | 9 | 7.5 | 7 | 7 | 9.5 | 2 | 8.5 | 100 |

**41% accuracy!**

# Highlights of $k$-Nearest Neighbor

1. Simple. ✓

2. No training.

3. Near optimal $E_{\text{out}}$. (N is sufficiently large)

} A **good!** method

4. Easy to justify classification to customer.

5. Can easily do multi-class.

6. Can easily adapt to regression or logistic regression

$$g(\mathbf{x}) = \frac{1}{k} \sum_{i=1}^{k} y_{[i]}(\mathbf{x})$$

$$g(\mathbf{x}) = \frac{1}{k} \sum_{i=1}^{k} [\![y_{[i]}(\mathbf{x}) = +1]\!]$$

7. **Computationally demanding.** ← we will address this next

# Thanks!