# 2/1/2024: Non-parametric Methods

| | Discriminative | Generative |
|---|---|---|
| **Parametric Methods**<br>- Fixed # of parameters to learn<br>- After learning, training data no longer needed | Logistic Regression<br>Softmax Regression<br>Directly model $p(y\|x)$<br><br>• Log. Reg. learn $w \in \mathbb{R}^d$<br>- Sof. Reg. learn $w^{(1)}, ..., w^{(c)} \in \mathbb{R}^d$ | Naive Bayes<br>model<br>$p(y)$ and $p(x\|y)$<br>↓      ↓<br>$\pi$    $\tilde{c}$<br>extracted from training data by counting |
| **Non-parametric Methods**<br>- Size of model proportionate to size of training dataset<br>- Usually because we store training dataset & use it to make prediction | K-Nearest Neighbors<br><br>Kernel Methods | |



← complex 1-NN decision boundary

**Logistic Regression:**
Can't fit data well (without adding more features)

**1-NN:** can fit training data always

## 1-Nearest Neighbor (1-NN)
Idea: Similar points usually have the same label

① Training step: Store training data in memory

② Test time: Given $x$ find most similar training example:
$$i^* = \operatorname*{argmin}_{i=1,...,n} \text{distance}(x, x^{(i)})$$

return $y^{(i^*)}$    (label of the most similar point)

Common distance is Euclidean distance
ie. $\|x - x^{(i)}\|$

K- Nearest Neighbors:
  - Find k closest training examples to test input x
  - Return most common label among those k
Why? Reduces effect of anomalous training examples

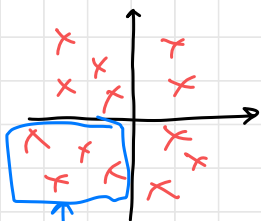## Pitfalls of k-NN

- Bias vs. Variance → Error in estimating best possible
  model in model family
  caused by overfitting

  Error b/c
  assumptions of model
  are wrong

  Very low b/c
  can represent
  any function

  Can be very large

- Curse of Dimensionality
  In high dimensions, you rarely have close neighbors



In $\mathbb{R}^2$, ~1/4 of points
are in same quadrant
as you

→ If in $\mathbb{R}^{1000}$
then only $\frac{1}{2^{1000}}$ points
are in same quadrant
Closest neighbor is still not that similar,
so they might not have same label

| K-NN | Logistic Regression |

K-NN
- Idea: Similar points
  have similar labels
- No good way to
  "regularize"
- No parameters we
  could tweak

Logistic Regression
- Only learns a linear
  decision boundary
- Learn parameters
  from data
- Regularization ($L_2$)

## Kernel Methods  Combine ideas from K-NN and Logistic Regression

Make a prediction on test example $x$ based on:
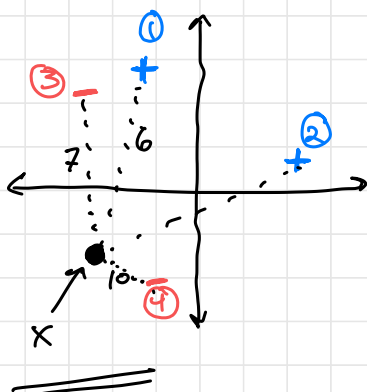
$$\sum_{i=1}^{n} \alpha_i \, K(x, x^{(i)})$$

total of
$n$
$\alpha_i$'s to $\leftarrow$ parameter
learn    to learn

$\in \mathbb{R}$,

"Kernel function"
measures similarity between 2 points

For binary classification:
- Logistic Regression: If $w^T x > 0$, predict $y = +1$
  If $w^T x < 0$,        $y = -1$

- Kernel-based classifier: If $\sum_{i=1}^{n} \alpha_i K(x, x^{(i)}) > 0$, predict $+1$
  $< 0$, predict $-1$

suppose: $K(x, x^{(1)}) = 6$
$x^{(2)}) = 1$    AND
$x^{(3)}) = 7$
$x^{(4)}) = 10$

$\alpha_1 = 1$
$\alpha_2 = 1$
$\alpha_3 = -1$
$\alpha_4 = -1$

Here: Score $= 6 \cdot 1 + 1 \cdot 1 + \boxed{7 \cdot -1} + \boxed{10 \cdot -1}$
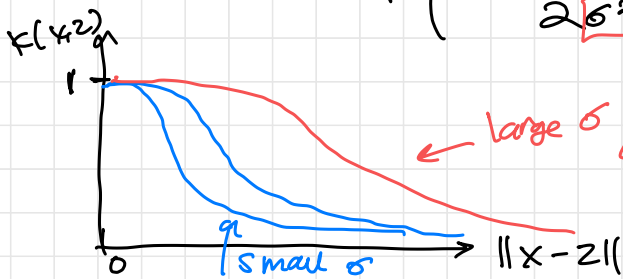
$= -11$   predict $-1$

One popular option for Kernel:
Radial Basis Function (RBF) Kernel

$$K(x, z) = \exp\left(-\frac{\|x - z\|^2}{2\sigma^2}\right)$$

hyperparameter
called "bandwidth"

large $\sigma$ = width of curve is larger
$\Longrightarrow$ points further away still
considered somewhat
similar

$K(x, z)$

large $\sigma$
small $\sigma$

$\|x - z\|$

How to learn $\alpha_i$'s?

Caveat: This is not recommended practice, but shows connection to logistic regression

Logistic Regression is a kernel method
using the kernel $\underline{K(x,z) = x^T z}$

## Logistic Regression

To predict on input $x$:
Compute $w^T x$
$= \left( \sum_{i=1}^{n} \alpha_i x^{(i)} \right)^T x = \sum_{i=1}^{n} \alpha_i x^{(i)T} x$

$= \boxed{\sum_{i=1}^{n} \alpha_i K(x^{(i)}, x)}$  Log. Reg. is a kernel method where $K(x,z) = x^T z$

Training: G.D.
$w^{(0)} \leftarrow \underline{0}$

$w^{(t)} \leftarrow w^{(t-1)}$
$+ \eta \cdot \frac{1}{n} \cdot \sum_{i=1}^{n} \underbrace{\sigma(-y^{(i)} w^{(t-1)T} x^{(i)}) \cdot y^{(i)}}_{\text{scalar}} \cdot x^{(i)}$

Key observation: update to $w$
is always $c_1 x^{(1)} + c_2 x^{(2)} + \cdots + c_n x^{(n)}$
So: Final $w$ can be written as
$\boxed{w = \sum_{i=1}^{n} \alpha_i x^{(i)}}$   $\alpha_i$'s are the weights of linear combination of the $x^{(i)}$'s

## Kernelized LR

Goal: learn $\alpha_i$'s for any given kernel function $K$