# Lecture 18: Non-parametric learning with K-nn

### K-Nearest Neighbor (K-nn) method

1. K-nn density estimation --- estimate the prior and class models

2. K-nn classification --- estimate the posterior directly

3. K-nn decision rule

4. k-nn error analysis

# Summary

Last lecture has introduced the non-parametric methods for density estimation. Given $n$ samples $\{x_1, x_2, ..., x_n\}$, the density at point x is estimated as,

$$p_n(x) = \frac{k_n/n}{V_n} = \frac{1}{n} \sum_{i=1}^{n} \delta_n(x - x_i)$$

There are two methods for designing the window sequence $V_1, V_2, ..., V_n$ at point $x$:

- the Parzen windows,
- the $k_n$-Nearest-Neighbor.

The Parzen window method is an unbiased estimate of the underlying density p(x)

$$\lim_{n \to \infty} p_n(x) = p(x), \quad \lim_{n \to \infty} \sigma^2(p_n(x)) = 0.$$

# Parzen window vs K-NN--- a comparison

Given n samples, in a non-parametric method, we estimate the density at x by

$$p_n(x) = \frac{k_n / n}{V_n}$$

Design I: Parzen window.

Fix the window size $V_n = \frac{V_1}{\sqrt{n}}$ and $k_n = k_n(x)$ is a function of x.

$$p_n(x) = \frac{k_n / n}{V_n} = \frac{k_n(x)}{V_1 \sqrt{n}}$$

Design II: K-nn method.

Fix the number of sample inside window $K_n = \sqrt{n}$ and $V_n = V_n(x)$ --- a function of x.

$$p_n(x) = \frac{k_n / n}{V_n} = \frac{1}{V_n(x) \sqrt{n}}$$

---

# Parzen window vs K-NN—a comparison

| Parzen window | K-nn |
|---|---|
| $V_n = \frac{V_1}{\sqrt{n}}$ | $K_n = \sqrt{n}$ |
| $p_n(x) = \dfrac{k_n / n}{V_n} = \dfrac{k_n(x)}{V_1 \sqrt{n}}$ | $p_n(x) = \dfrac{k_n / n}{V_n} = \dfrac{1}{V_n(x) \sqrt{n}}$ |
| 1. Here the number of samples falling in a window can be counted explicitly | 1. The window size does not have an explicitly form. Very irregular plot, the diameter of the window is |
| $K_n(x) = \sum_{i=1}^{n} \varphi(\frac{x - x_i}{h_n})$ | $\phi(V_n(x)) = 2\,|x - x_K^*|, \quad x_K^* \in \{x_1,...,x_n\}$ |
| 2. We can prove the convergence | 2. The integration of $p_n(x)$ is infinity |
| | $\int p_n(x)\,dx = \infty$ |

# Examples of K-nn density
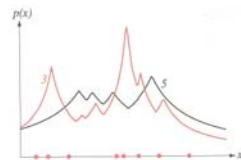
$$p_n(x) = \frac{k_n/n}{V_n} = \frac{1}{V_n(x)\sqrt{n}}$$



FIGURE 4.10. Eight points in one dimension and the $k$-nearest-neighbor density estimates, for $k = 3$ and 5. Note especially that the discontinuities in the slopes in the estimates generally lie *away* from the positions of the prototype points.
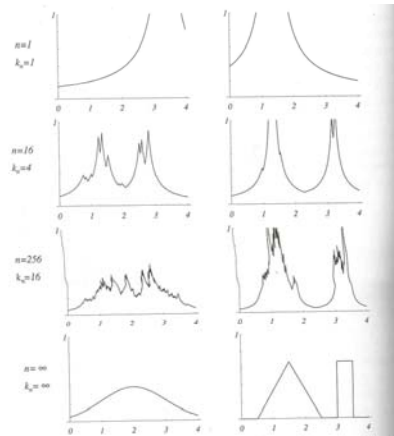
$$p_1(x) = \frac{1}{2|x - x_1|}$$

FIGURE 4.12. Several $k$-nearest-neighbor estimates of two unidimensional densities: Gaussian and a bimodal distribution. Notice how the finite $n$ estimates can be quite "spiky."

---

# Non-parametric Classification

Given $n$ labeled samples $\chi = \{x_1, x_2, ..., x_n\}$, and we have

| classes | $\omega_1, \omega_2, ..., \omega_c$ | |
|---|---|---|
| #samples | $n_1, n_2, ..., n_c$ | $\sum_{i=1}^{c} n_i = n.$ |
| $p(\omega_i|\chi)$ | $\frac{n_1}{n}, \frac{n_2}{n}, ..., \frac{n_c}{n},$ | the prior probabilities. |

At each point $x$, a cell of volume $V_n$ covers total $k$ samples, with $k_i$ samples in each class $\omega_i$, $i = 1, 2, ..., c$. The likelihood is estimated as

$$p(x|\omega_i, \chi) = \frac{k_i/n_i}{V_n}, \qquad i = 1, 2, ..., c.$$

Note that the above class models are estimated by the Parzen windows – all $c$ classes share the same window $V_n$ at $x$, and it is possible that $k_i = 0$ for some classes. In general, One may use different windows $V_n$ for different classes.

# Non-parametric Classification

Given the learned prior model and class models, we compute the posterior probability,

$$P(\omega_i|\mathbf{x}, \chi) = \frac{p(\mathbf{x}|\omega_i, \chi)p(\omega_i|\chi)}{\Sigma_{j=1}^c p(\mathbf{x}|\omega_j, \chi)p(\omega_j|\chi)} = \frac{k_i}{k}.$$

– it is the fraction of the samples within the cell that are labeled $\omega_i$.

The Bayes rule chooses

$$\omega^*(\mathbf{x}) = \arg\max_i\{k_1, k_2, ..., k_c\}.$$

It turns out that the K-nn classification is also a special case of Bayes decision !

# Nearest Neighbor Decision Rule

The Bayes decision rule in the non-parametric case indicates that we can go directly for a decision, and thus bypass the steps of estimating $p(\omega_i|\chi)$ and $p(\mathbf{x}|\omega_i, \chi)$. This leads to two methods:

- the nearest neighbor decision rule,
- the $k_n$-nearest neighbor decision rule.

Intuitively, both methods partition the feature space into $c$ disjoint subspaces,

$$\Omega = \cup_{i=1}^n \Omega_i, \quad \Omega_i \cap \Omega_j = \emptyset, i \neq j.$$
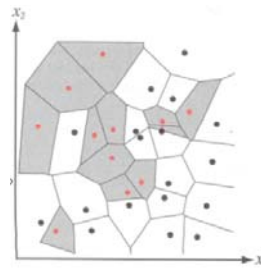
It is possible that $\Omega_i$ is not connected for some $i$.

# Nearest Neighbor Decision Rule

Let $\{(x_1, \omega(x_i)), ..., (x_n, \omega(x)_n)\}$ be the labeled samples, for any $x \in \Omega$,

$$\omega_{NN}(x) = \omega(x^*), \quad x^* = \arg \min_{x_j}\{\|x - x_j\|, \; j = 1, 2, ..., n\}$$

Remark I: this simply partitions the feature space by a Voronoi diagram, with each sample $x_i$ occupies a cell.

# Nearest Neighbor Decision Rule

**Asymptotic Behavior of the NN decision rule.**

When large samples $\chi = \{x_1, x_2, ..., x_n\}$ are available, for any point $x$, its nearest sample $x^*$ in $\chi$ *approaches* $x$, therefore

$$p(\omega_i|x^*) \approx p(\omega_i|x).$$

Therefore, the chance of $\omega(x^*) = \omega_i$ is simply $p(\omega_i|x)$, that is,

$$\omega_{NN}(x) = \omega(x^*) \sim p(\omega|x)$$

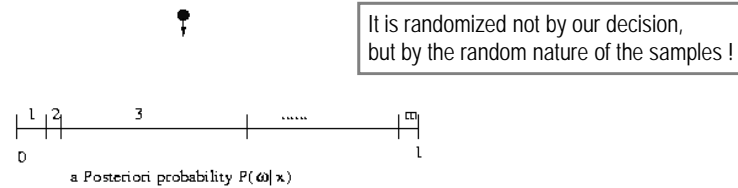## Randomized decision

**Asymptotic Behavior of the NN decision rule.**

The Nearest-Neighbor decision rule is a randomized decision procedure, draw a random number $s \in [0,1]$, and decide $\omega_{NN}(x)$ according to the posterior probability $p(\omega|x)$.

> It is randomized not by our decision, but by the random nature of the samples !

a Posteriori probability $P(\omega|x)$

If $n$ is large enough, $P(\omega|x)$ is close to the true posterior probability.

In fact, we are assuming that the samples are drawn from the true posterior probability.

## Two Special Cases

Recall that we have discussed this randomized decision before in comparison with the Bayes decision rule. We have proven that the randomized decision is not as good as the Bayes decision rule.

Two special cases,

1. when $P(\omega_i|x) \approx 1.0$ for a one class, then in most of the time, $\omega_{NN}(x) = \omega_{Bayes}(x)$.

2. when $P(\omega_i|x) \approx \frac{1}{c}$, i.e. a uniform distribution, then $\omega_{NN}(x) = \omega_{Bayes}(x)$ with one out of $c$ times.

In both cases, the performance of the Nearest-Neighbor decision rule matches that of the Bayes rule.

Now we analyze the error rate of the NN rule qumtitatively.

# NN error analysis

**Error rate of the NN-decision rule at x with $n$ samples.**

$$P_n(e|\mathbf{x}) = \int p_n(e, \mathbf{x}^*|\mathbf{x})d\mathbf{x}^* = \int P_n(e|\mathbf{x}^*, \mathbf{x})p(\mathbf{x}^*|\mathbf{x})d\mathbf{x}^*$$

as we know,

$$p(\mathbf{x}^*|\mathbf{x}) = \delta(\mathbf{x} - \mathbf{x}^*) \quad \text{as} \quad n \to \infty,$$

where $\delta()$ is the Dirac delta function, and

$$P_n(e|\mathbf{x}^*, \mathbf{x}) = 1 - \sum_{i=1}^{c} P(\omega(\mathbf{x}) = \omega_i, \omega(\mathbf{x}^*) = \omega_i|\mathbf{x}^*, \mathbf{x})$$
$$= 1 - \sum_{i=1}^{c} P(\omega_i|\mathbf{x}^*)P(\omega_i|\mathbf{x})$$

---

# NN error analysis

**Error rate of the NN-decision rule at x with $n$ samples.**

$$\lim_{n \to \infty} P_n(e|\mathbf{x}) = \int [1 - \sum_{i=1}^{c} P(\omega_i|\mathbf{x}^*)P(\omega_i|\mathbf{x})]\delta(\mathbf{x} - \mathbf{x}^*)d\mathbf{x}^*$$
$$= 1 - \sum_{i=1}^{c} P^2(\omega_i|\mathbf{x})$$

and the expected error rate is

$$P = \lim_{n \to \infty} \int P_n(e|\mathbf{x})p(\mathbf{x})d\mathbf{x} = \int [1 - \sum_{i=1}^{c} P^2(\omega_i|\mathbf{x})]p(\mathbf{x})d\mathbf{x}$$

In comparison, the optimal Bayes rate is

$$P^* = \int [1 - max\{P(\omega_{Bayes}(\mathbf{x})|\mathbf{x})\}]p(\mathbf{x})d\mathbf{x}$$

The Bayes rate is obtained by assuming we know the true densities.

# NN error analysis

**The error bound of the NN-decision rule:**

$$P^* \leq P \leq P^*(2 - \frac{c}{c-1}P^*)$$

[Proof] The lower bound is obvious, we only prove the upper bound.
Let $\omega_m = \omega_{Bayes}(\mathbf{x})$, and $P^*(e|\mathbf{x}) = 1 - P(\omega_m|\mathbf{x})$, then we have,

$$\sum_{i=1}^{c} P^2(\omega_i|\mathbf{x}) = P^2(\omega_m|\mathbf{x}) + \sum_{i \neq m} P^2(\omega_i|\mathbf{x})$$

$$\geq (1 - P^*(e|\mathbf{x}))^2 + \frac{P^{*2}(e|\mathbf{x})}{c-1}.$$

so

$$1 - \sum_{i=1}^{c} P^2(\omega_i|\mathbf{x}) \leq P^*(e|\mathbf{x})(2 - \frac{c}{c-1}P^*(e|\mathbf{x}))$$
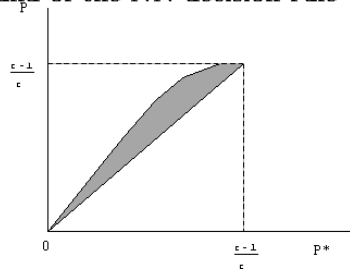
Take an integral on both sides, the conclusion follows.

A lemma: for $n$ positive real numbers $z_1, z_2, ..., z_n$ with constraint $\sum_{i=1}^{n} z_i = a$, then $\sum_{i=1}^{n} n z_i^2 \geq \frac{a^2}{n}$. Prove it by lagrange multipliers.

---

# NN error analysis

**The error bound of the NN-decision rule**



$P$ reaches $P^*$ at two extreme cases:

Case I: when no information is available, $P = P^* = \frac{c-1}{c}$.

Case II: when there is no uncertainty, $P = P^* = 0$.

But this error bound is obtained as $n \to \infty$, for finite samples the convergence rate can be arbitrarily slow. The error rate $P_n(\epsilon)$ could be very bad.
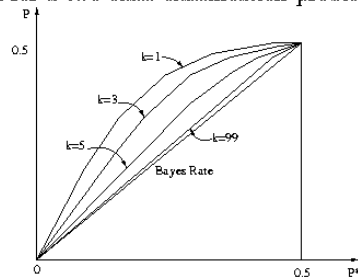
# Knn-Decision Rule

$k$-Nearest Neighbor Decision Rule:

At each point x, assign x the label that is most frequently represented among the $k$-nearest samples.

The error bound for a two-class classification problem.



The bigger $k$ is, the more accurate the estimation of the class models. However, the larger $k$ is, the bigger the window is, thus the less accurate the estimation is. Given finite training samples, the optimal $k$ is a trade-off between these two aspects.

---

## Discussion on distance measure

The K-nn method is quite popular due to its simplicity in learning and classification
This lecture studies some in-depth issues related to K-nn.

1  Means for expedite the K-nn algorithm
      a).  Editing and pruning (redundant) training samples
      b).  Search tree
      c).  Computing partial distance (pruning samples over a certain bound value)

2.   Metircs and distance

3.   Invariance to transformations

4.  Tangent distance

# K-nn distance measure

For any two points x and y in a d-space

- Minkowski metric

$$D(x, y) = (\sum_{i=1}^{d} |x_i - y_i|^k)^{1/k}$$
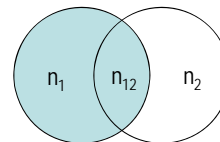
2. Manhattan metric

$$D(x, y) = \sum_{i=1}^{d} |x_i - y_i|$$

3. Mahalanobis metric

$$D(x, y) = X' \sum^{-1} Y$$

4. Tanimoto metric for two sets

$$D(s_1, s_2) = \frac{n_1 + n_2 - 2n_{12}}{n_1 + n_2 - n_{12}}$$
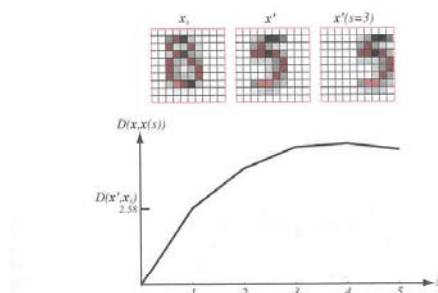
---

# Distance is sensitive to transforms

The distance from "5" in the middle is closer to "8" on the left, in comparison with the distance to and slightly shifted "5" on the right.



FIGURE 4.20. The uncritical use of Euclidean metric cannot address the problem of translation invariance. Pattern x' represents a handwritten 5, and x'(s = 3) represents the same shape but shifted three pixels to the right. The Euclidean distance $D(x', x'(s = 3))$ is much larger than $D(x', x_8)$, where $x_8$ represents the handwritten 8. Nearest-neighbor classification based on the Euclidean distance in this way leads to very large errors. Instead, we seek a distance measure that would be insensitive to such translations, or indeed other known invariances, such as scale or rotation.

# Invariant distance to transforms

There are two ways to make the nn-distance work under various transforms.

Method 1:

For exch training example $x_i$, we make many copies with transformations (translation, rotation and scaling).

Method 2:

We compute a distance on-line which is the minimum distance between y and a transformed x, by a transform function with parameter a in a transformation group G. For example, a is a vector for translations, rotation, scaling etc.
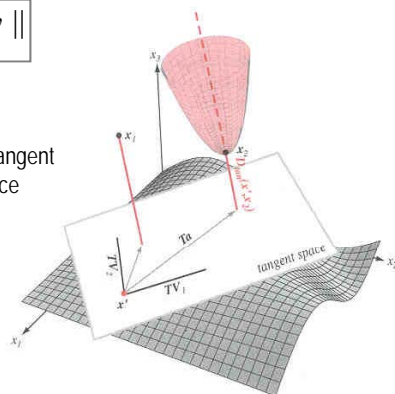
$$D(x, y) = \min_{a \in G} \| f(x; a) - y \|$$

# Tangent distance

One may approximate the transform function by a first order Taylor expansion, thus we have a tangent distance

$$D(x, y) = \min_{a \in G} \| (x + Ta) - y \|$$

T is a matrix for the gradients. X+Ta spans a tangent plane at x. D(x,y) is the perpendicular distance from y to the tangent plane.

# Example of the tangent plane