

MLISP: Machine Learning in Signal Processing

Lecture 4

Prof. V. I. Morgenshtern

Scribe: M. Solomon

Illustrations: The elements of statistical learning, Hastie, Tibshirani, Friedman

Agenda:

1. Basics of statistical decision theory
2. Bayesian optimality
3. Conditional expectation as the optimal predictor
4. Conditional expectation in the linear case
5. Maximum conditional probability as the optimal classifier

1 Basics of statistical decision theory

In this lecture we will touch upon the basics of statistical decision theory. Bayesian setup is the case in which the joint probability distribution from which our data was generated is known to us. We will see that in this case, there is an optimal solution to the regression and the classification problems. These are known as ‘Bayes optimal predictors/classifiers’. In practice, the joint probability distribution is never known. We will see how the two algorithms, the nearest neighbors predictor and the linear model, attempt to approximate the unachievable optimal Bayesian solution.

2 Bayes optimal predictor: condition expectation

Let $\mathbf{x} \in \mathbb{R}^p$ denote a real valued random input vector, and $y \in \mathbb{R}$ a real valued output vector. Assume that (\mathbf{x}, y) pairs are drawn from a joint probability distribution $\mathbb{P}_{\mathbf{x},y}(\cdot, \cdot)$.

Notation: we will use the notation and write \mathbf{x} for random vectors, y for random scalars (not as bold as vectors), \mathbf{x} for non-random vectors, y for non-random scalars.

Example: For example we might have the joint distribution $\mathbb{P}_{\mathbf{x},y}(\cdot, \cdot)$ defined (implicitly) as follows. The house size, $x_1 \sim \mathcal{U}[70, 200]$ sq. meters. The number of bedrooms, $x_2 = x_1/35 + \mathcal{U}[-1, 0, 1]$. The house price $y = \alpha_1 x_1 + \alpha_2 x_2 + \mathcal{N}(0, 100000)$ US dollars. Here and in the future, \mathcal{U} denotes the uniform distribution and \mathcal{N} denotes the Gaussian distribution.

Recall, we seek a function $h(\cdot)$ for predicting y given \mathbf{x} . Let us use the squared error loss as before:

$$L(y, h(\mathbf{x})) = (y - h(\mathbf{x}))^2. \quad (1)$$

Then the expected prediction error is given by:

$$\begin{aligned} \text{EPE}(h) &= \mathbb{E} \{ (y - h(\mathbf{x}))^2 \} \\ &= \int (y - h(\mathbf{x}))^2 p(\mathbf{x}, y) d\mathbf{x} dy \end{aligned} \quad (2)$$

where $p(\cdot, \cdot) = p_{\mathbf{x}, y}(\cdot, \cdot)$ denotes the joint probability density function of (\mathbf{x}, y) . In the Bayesian formulation, the goal is to find $h(\cdot)$ that minimizes $\text{EPE}(h)$.

Using the formula for conditional expectation:

$$\text{EPE}(h) = \mathbb{E}_{\mathbf{x}} \mathbb{E}_{y|\mathbf{x}} \{ (y - h(\mathbf{x}))^2 | \mathbf{x} \}. \quad (3)$$

Therefore, we can minimize EPE point-wise:

$$\hat{h}(\mathbf{x}) = \arg \min_c \mathbb{E}_{y|\mathbf{x}} \{ (y - c)^2 | \mathbf{x} = \mathbf{x} \} \quad (4)$$

for each \mathbf{x} . The solution is $\hat{h}(\mathbf{x}) = \mathbb{E} \{ y | \mathbf{x} = \mathbf{x} \}$, which is also known as the *regression function*.

Proof: The minimum is achieved at the point where the derivative of the function is zero:

$$\frac{d}{dc} \mathbb{E}_{y|\mathbf{x}} \{ (y - c)^2 | \mathbf{x} = \mathbf{x} \} = 0 \quad (5)$$

$$\Leftrightarrow - \int p_{y|\mathbf{x}}(y | \mathbf{x} = \mathbf{x}) 2(y - c) dy = 0 \quad (6)$$

$$\Leftrightarrow c = \int_y p_{y|\mathbf{x}}(y | \mathbf{x} = \mathbf{x}) y dy \quad (7)$$

$$\Leftrightarrow c = \mathbb{E} \{ y | \mathbf{x} = \mathbf{x} \}. \quad (8)$$

2.1 Nearest neighbors as an approximation to optimal Bayes solution

In practice we cannot use the $\hat{h}(\mathbf{x}) = \mathbb{E} \{ y | \mathbf{x} = \mathbf{x} \}$ formula directly because we don't know the joint probability density $p_{\mathbf{x}, y}(\cdot, \cdot)$, which is the whole point of learning.

The nearest-neighbor methods attempt to directly implement the optimal recipe using the training data $(\mathbf{x}^{(i)}, y^{(i)})$ to estimate $\mathbb{E} \{ y | \mathbf{x} = \mathbf{x} \}$. Specifically, at each point \mathbf{x} , we average all those $y^{(i)}$'s with the corresponding inputs $\mathbf{x}^{(i)} \approx \mathbf{x}$:

$$\hat{h}_{\text{NN}}(\mathbf{x}) = \text{Ave} \left(y^{(i)} | \mathbf{x}^{(i)} \in N_k(\mathbf{x}) \right) = \frac{1}{|N_k(\mathbf{x})|} \sum_{i: \mathbf{x}^{(i)} \in N_k(\mathbf{x})} y^{(i)} \quad (9)$$

where ‘‘Ave’’ denotes the average, and $N_k(\mathbf{x})$ is the neighborhood consisting of the k points in the training set that are the closest to \mathbf{x} .

If we have a lot of data: $n, k \rightarrow \infty$ such that $k/n \rightarrow 0$, then: $\hat{h}_{\text{NN}}(\mathbf{x}) \rightarrow \mathbb{E} \{ y | \mathbf{x} = \mathbf{x} \} = \hat{h}(\mathbf{x})$.

To see this intuitively, consider a special case and assume that we only have two features, $\mathbf{x} \in [0, 1] \times [0, 1]$ and that $p_{\mathbf{x}}(\mathbf{x})$ is uniform over the set $\mathcal{X} = [0, 1] \times [0, 1]$. Let $|\mathcal{X}|$ denote the area of \mathcal{X} , $|\mathcal{X}| = 1$ in this case. By definition,

$$\mathbb{E}_{y|\mathbf{x}} \{ y | \mathbf{x} = \mathbf{x} \} = \int_y y p_{y|\mathbf{x}}(y | \mathbf{x} = \mathbf{x}) dy. \quad (10)$$

Let $N(\mathbf{x})$ be the area around \mathbf{x} that is roughly equal to $N_k(\mathbf{x})$. Assume for simplicity that the neighborhoods $N_k(\mathbf{x})$ are all the same size for all \mathbf{x} (follows from uniformity assumption). Then, starting from (10):

$$\begin{aligned}\mathbb{E}_{\mathbf{y}|\mathbf{x}}\{\mathbf{y}|\mathbf{x}=\mathbf{x}\} &= \frac{1}{|N(\mathbf{x})|} \int_{\mathbf{x}' \in N(\mathbf{x})} \int_{\mathbf{y}} y p_{\mathbf{y}|\mathbf{x}}(y|\mathbf{x}=\mathbf{x}) p_{\mathbf{x}}(\mathbf{x}') d\mathbf{x}' dy \\ &\approx \frac{1}{|N(\mathbf{x})|} \int_{\mathbf{x}' \in N(\mathbf{x})} \int_{\mathbf{y}} y p_{\mathbf{y}|\mathbf{x}}(y|\mathbf{x}=\mathbf{x}') p_{\mathbf{x}}(\mathbf{x}') d\mathbf{x}' dy \\ &= \frac{1}{|N(\mathbf{x})|} \int_{\mathbf{x}' \in N(\mathbf{x})} \int_{\mathbf{y}} y p_{\mathbf{y},\mathbf{x}}(\mathbf{x}', y) d\mathbf{x}' dy \\ &\approx \text{Ave}\left(\mathbf{y}^{(i)}|\mathbf{x}^{(i)} \in N_k(\mathbf{x})\right).\end{aligned}$$

Above, in the second equality we assumed that the function $p_{\mathbf{y}|\mathbf{x}}(y|\mathbf{x}=\mathbf{x}')$ is a smooth function of \mathbf{x}' and does not change much in the neighborhood $\mathbf{x}' \in N(\mathbf{x})$.

Isn't this an ideal universal learning method?

In low dimensions this is true. However when dimensionality gets large we run into very serious problems. As we will see soon, the metric size of k -neighborhood gets very large as the dimension gets large. Therefore, nearest neighborhood becomes a very poor surrogate for conditioning.

2.2 Linear regression as an approximation to Bayes solution in the linear setting

How does the least squares algorithm we discussed before fit into the Bayesian theory? Suppose that from our prior knowledge of the problem domain we decide to make an additional assumption on the functional form of the regression function $h(\cdot)$. Specifically, assume that $h(\cdot)$ is linear: $h(\mathbf{x}) = \mathbf{x}^T \boldsymbol{\theta}$.

Plugging this into (2) and differentiating to find the minimum of $\text{EPE}(\boldsymbol{\theta})$, we can solve for $\boldsymbol{\theta}$:

$$\boldsymbol{\theta}^* = \left[\mathbb{E}\{\mathbf{xx}^T\}\right]^{-1} \mathbb{E}\{\mathbf{xy}\}. \quad (11)$$

Proof:

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \mathbb{E}\{[\mathbf{y} - \mathbf{x}^T \boldsymbol{\theta}]^2\} \quad (12)$$

$$\Leftrightarrow \frac{d}{d\hat{\boldsymbol{\theta}}} \int p(\mathbf{x}, y) [y - \mathbf{x}^T \hat{\boldsymbol{\theta}}]^2 d\mathbf{x} dy = 0 \quad (13)$$

$$\Leftrightarrow \int p(\mathbf{x}, y) 2\mathbf{x} [y - \mathbf{x}^T \hat{\boldsymbol{\theta}}] d\mathbf{x} dy = 0 \quad (14)$$

$$\Leftrightarrow \int p(\mathbf{x}, y) \mathbf{xx}^T \hat{\boldsymbol{\theta}} d\mathbf{x} dy = \int p(\mathbf{x}, y) \mathbf{xy} d\mathbf{x} dy \quad (15)$$

$$\Leftrightarrow \mathbb{E}[\mathbf{xx}^T] \hat{\boldsymbol{\theta}} = \mathbb{E}[\mathbf{xy}] \quad (16)$$

$$\Rightarrow \hat{\boldsymbol{\theta}} = \left[\mathbb{E}\{\mathbf{xx}^T\}\right]^{-1} \mathbb{E}\{\mathbf{xy}\}. \quad (17)$$

Let's compare this to least-squares solution we derived in the previous lecture:

$$\hat{\boldsymbol{\theta}}_{\text{LS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (18)$$

where \mathbf{X} is the data matrix:

$$\mathbf{X} = \begin{bmatrix} - & (\mathbf{x}^{(1)})^\top & - \\ - & (\mathbf{x}^{(2)})^\top & - \\ - & \vdots & - \\ - & (\mathbf{x}^{(n)})^\top & - \end{bmatrix}.$$

To see the relationship between the two formulas, consider the j th column of $\mathbf{X}^\top \mathbf{X}$:

$$[\mathbf{X}^\top \mathbf{X}]_j = \sum_{i=1}^n \mathbf{x}^{(i)} x_j^{(i)} \quad (19)$$

where $x_j^{(i)}$ is the j th variable in the i th example in the training set and $\mathbf{x}^{(i)} = [x_0^{(i)}, \dots, x_{p-1}^{(i)}]^\top$. By the law of large numbers this converges (when the number of training examples, n becomes large) to $n \mathbb{E}\{\mathbf{x}x_j\}$, which is n times the j th column of $\mathbb{E}\{\mathbf{x}\mathbf{x}^\top\}$.

Similarly, note that

$$\mathbf{X}^\top \mathbf{y} = \sum_{i=1}^n \mathbf{x}^{(i)} y^{(i)}. \quad (20)$$

By the law of large numbers this converges (when the number of training examples, n becomes large) to $n \mathbb{E}\{\mathbf{x}y\}$. We conclude that

$$\hat{\boldsymbol{\theta}}_{\text{LS}} \approx \left[n \mathbb{E}\{\mathbf{x}\mathbf{x}^\top\} \right]^{-1} n \mathbb{E}\{\mathbf{x}y\} = \hat{\boldsymbol{\theta}}. \quad (21)$$

2.3 Summary

Both k -nearest neighbors and least squares end up approximating conditional expectation by average. But they differ dramatically in terms of model assumptions:

- Least Squares assumes that $h(\mathbf{x})$ is well approximated by a global linear function.
- k -nearest neighbors assumes that $h(\mathbf{x})$ is well approximated by a local constant function.

3 Bayes optimal classifier: maximum conditional probability

So far in this lecture we have considered a regression problem, let's return to the classification problem from the previous lecture:

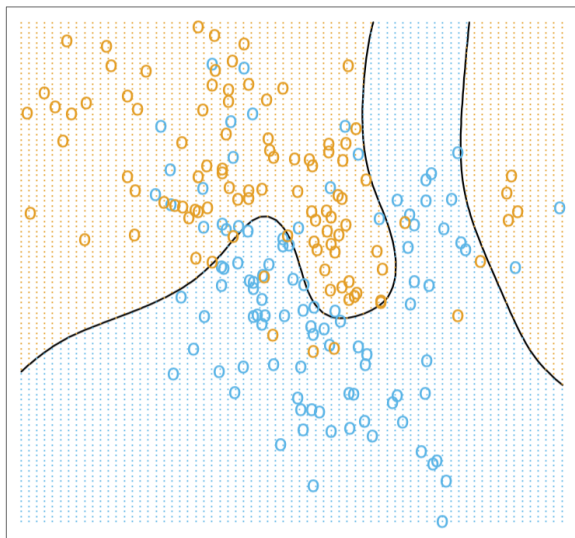


Figure 1: Bayes optimal classifier

We did not specify before how the data points in this example were generated. Now we reveal the secret:

- Generate 10 means \mathbf{m}_k from $\mathcal{N}([0, 1]^\top, \mathbf{I})$ and label this class BLUE. Generate 10 means \mathbf{m}_k from $\mathcal{N}([1, 0]^\top, \mathbf{I})$ and label this class ORANGE. These 20 points are fixed once and for all and are assumed nonrandom and known.
- For each class (BLUE and ORANGE), generate 100 observations as follows: pick an m_k at random with probability of 1/10, then generate points according to $\mathcal{N}(\mathbf{m}_k, \mathbf{I}/5)$.

This is a Bayesian problem because we have specified the data-generating distribution precisely. Therefore, we can calculate the optimal solution by minimizing the expected prediction error:

$$EPE(h) = \mathbb{E}(L(\mathbf{y}, h(\mathbf{x}))) \quad (22)$$

where we use the misclassification loss function that is defined as follows:

$$L(\mathbf{y}, h(\mathbf{x})) = \begin{cases} 0, & \text{iff } h(\mathbf{x}) = \mathbf{y} \\ 1, & \text{iff } h(\mathbf{x}) \neq \mathbf{y}. \end{cases} \quad (23)$$

Similarly to the calculation we did for regression,

$$h(\mathbf{x}) = \arg \min_{c \in \{0,1\}} \mathbb{E}_{y|\mathbf{x}}(L(y, c)|\mathbf{x} = \mathbf{x}) \quad (24)$$

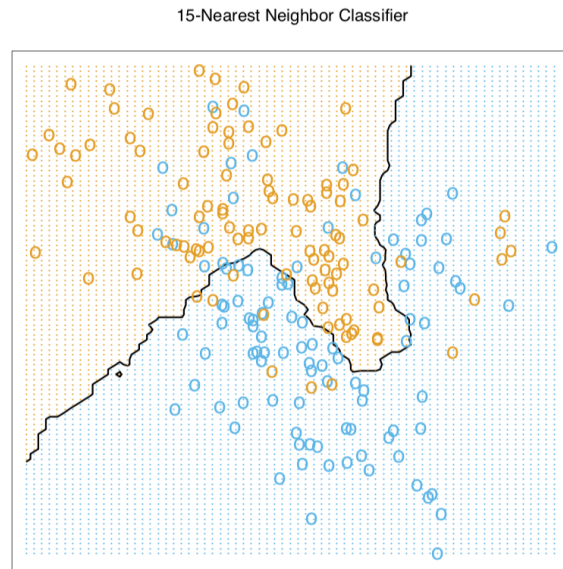
$$= \arg \min_{c \in \{0,1\}} \mathbb{P}(y = 0|\mathbf{x} = \mathbf{x}) \cdot L(0, c) + \mathbb{P}(y = 1|\mathbf{x} = \mathbf{x}) \cdot L(1, c) \quad (25)$$

$$= \arg \min_{c \in \{0,1\}} \mathbb{P}(y \neq c|\mathbf{x} = \mathbf{x}) \quad (26)$$

$$= \arg \min_{c \in \{0,1\}} 1 - \mathbb{P}(y = c|\mathbf{x} = \mathbf{x}) \quad (27)$$

$$= \arg \max_{c \in \{0,1\}} \mathbb{P}(y = c|\mathbf{x} = \mathbf{x}). \quad (28)$$

Therefore, the optimal algorithm is simply to choose the class that has the largest conditional probability given the data. This is called the *Bayes rule* for classification. The Bayes-optimal decision boundary is shown in Figure 1. It is instructive to compare this decision boundary to the one produced by the 15-nearest neighbors classifier:



We observe that the decision boundaries produced by the two algorithms are very similar!