

Regression (Theory!)

Definition 3.1 (Regression) Let \mathbf{X} be a D -dimensional input and \mathbf{Y} be a K -dimensional quantitative output with joint density $\rho(\mathbf{x}, \mathbf{y})$. In **regression** we seek, given a training set $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$, for a function $\mathbf{f} : \mathbb{R}^D \rightarrow \mathbb{R}^K$ that minimizes the error between \mathbf{Y} and its prediction $\hat{\mathbf{Y}} = \mathbf{f}(\mathbf{X})$.

Definition 3.4 (Expected (squared) prediction error in regression) Let \mathbf{X}, \mathbf{Y} be input and quantitative output with joint density $\rho(\mathbf{x}, \mathbf{y})$ and $\mathbf{f} : \mathbb{R}^D \rightarrow \mathbb{R}^K$ a function to predict \mathbf{Y} from \mathbf{X} . For the squared loss, the **expected (squared) prediction error** or **risk** of \mathbf{f} , $\text{EPE}(\mathbf{f})$, is defined as

$$\begin{aligned}\text{EPE}(\mathbf{f}) &= \mathbb{E}(L_2(\mathbf{Y}, \mathbf{f}(\mathbf{X}))) \\ &= \int_{\mathbb{R}^K} \int_{\mathbb{R}^D} L_2(\mathbf{y}, \mathbf{f}(\mathbf{x})) \rho(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y}.\end{aligned}$$

Theorem 3.1 Let \mathbf{X}, \mathbf{Y} be input and quantitative output with joint density $\rho(\mathbf{x}, \mathbf{y})$. The function $\mathbf{f} : \mathbb{R}^D \rightarrow \mathbb{R}^K$ that minimizes (under appropriate conditions) the expected (squared) prediction error $\text{EPE}(\mathbf{f})$ is given by

$$\mathbf{f}(\mathbf{x}) = \arg \min_{\mathbf{g}(\mathbf{x})} \text{EPE}(\mathbf{g}) = \mathbb{E}(\mathbf{Y} | \mathbf{X} = \mathbf{x})$$

$\mathbb{E}(\mathbf{Y} | \mathbf{X} = \mathbf{x})$ is often called **regressor** / **regression function**.

$$\left. \begin{array}{l} X: \Omega \rightarrow \mathbb{R}^D \\ Y: \Omega \rightarrow \mathbb{R} \end{array} \right\} \begin{array}{l} p(x, y) = ? \\ \Rightarrow p(y|x) = ? \end{array}$$

$$\begin{array}{l} \text{Idea: } Y \approx \hat{Y} = f(X) \\ f(X) = ? \end{array}$$

$$\begin{array}{l} \text{"prediction error of } f\text{"} \\ \text{EPE}(f) \end{array}$$

Warning! Did not use $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$!

Which is "best" f ?

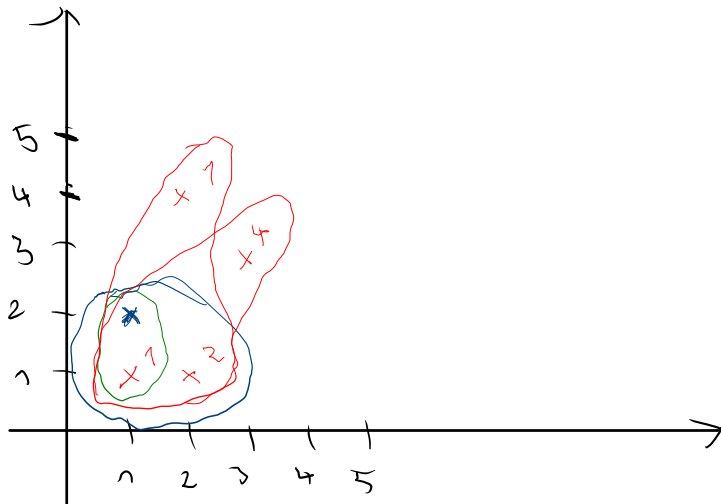
$$f(x) = \mathbb{E}(Y | X=x)$$

2 kNN regression

We are given the training data

$$\mathcal{T} = \{((1, 1)^\top, 1), ((2, 1)^\top, 2), ((3, 3)^\top, 4), ((2, 4)^\top, 1)\}.$$

Use kNN regression with $k = 1, 2, 3$ to do a prediction for $\mathbf{x} = (1, 2)$.



$$k = 1 \Rightarrow f((1, 2)^\top) = 1$$

$$k = 2 \Rightarrow f((1, 2)^\top) = 1.5$$

$$k = 3 \Rightarrow \begin{aligned} f((1, 2)^\top) &= 1.\bar{3} \\ f((1, 2)^\top) &= 2\frac{1}{3} \end{aligned}$$

Classification

Definition 3.2 (Classification) Let \mathbf{X} be a D -dimensional input and G be a one-dimensional qualitative output with density $\rho(\mathbf{x}, g)$ and range R_G . In **classification** we seek, given a training set $\{(\mathbf{x}_i, g_i)\}_{i=1}^N$, for a function $f: \mathbb{R}^D \rightarrow R_G$ that minimizes the error between G and its prediction $\hat{G} = f(\mathbf{X})$.

Definition 3.5 (Expected 0-1 prediction error in classification) Let \mathbf{X} , G be input and a one-dimensional output with mixed joint density $\rho(\mathbf{x}, g)$ and $f: \mathbb{R}^D \rightarrow R_G$ a function to predict G from \mathbf{X} . For the 0-1 loss, the **expected prediction error of f** , $\text{EPE}(f)$, is defined as

$$\begin{aligned} \text{EPE}(f) = \mathbb{E}(L_{0-1}(G, f(\mathbf{X}))) &= \sum_{g \in R_G} \int_{\mathbb{R}^D} L_{0-1}(g, f(\mathbf{x})) \rho(\mathbf{x}|g) d\mathbf{x} p_G(g) \\ &= \int_{\mathbb{R}^D} \sum_{g \in R_G} (L_{0-1}(g, f(\mathbf{x})) p(g|\mathbf{x})) \rho_{\mathbf{X}}(\mathbf{x}) d\mathbf{x}. \end{aligned}$$

Theorem 3.2 Let \mathbf{X} be a continuous D -dimensional input and G be the qualitative discrete output ($R_G = \{1, \dots, r\}$) with mixed joint density $\rho(\mathbf{x}, g)$. The function $f: \mathbb{R}^D \rightarrow R_G$ that minimizes (under appropriate conditions) the expected prediction error $\text{EPE}(f)$ with respect to the 0-1 loss is given by

$$f(\mathbf{x}) = \arg \min_{g \in R_G} [1 - p(g|\mathbf{x})],$$

which can be simplified to

$$f(\mathbf{x}) = \mathbf{g} \quad \text{if} \quad p(\mathbf{g}|\mathbf{x}) = \max_{g \in R_G} p(g|\mathbf{x})$$

This minimizer is the **Bayes classifier**.

$$\left. \begin{aligned} \mathbf{X}: \Omega &\rightarrow \mathbb{R}^D \\ G: \Omega &\rightarrow \{1, \dots, r\} \end{aligned} \right\} \begin{aligned} &\rho(\mathbf{x}, g) = ? \\ &\Rightarrow \rho(g|\mathbf{x}) = ? \end{aligned}$$

$$\text{Idea: } G \approx \hat{G} = f(\mathbf{X})$$

$$f(\mathbf{x}) = ?$$

"prediction error of f ":

$$\text{EPE}(f) \leftarrow 0-1 \text{ loss}$$

Which is best f ?

$$f(\mathbf{x}) = \arg \max_{g \in \{1, \dots, r\}} p(g|\mathbf{x})$$

\Rightarrow construction of classifiers

classes: 1, 2, 3

$$p(1|\mathbf{x}) = 0.1$$

$$p(\textcircled{2}|\mathbf{x}) = 0.5$$

$$p(3|\mathbf{x}) = 0.4$$

3 kNN classification

Previously, we studied the the well-known [Iris data set](#) and made it available in a Jupyter notebook. Moreover, we had an example Jupyter notebook in which we performed classification by the kNN classification method.

Combine both notebooks such that you carry out kNN classification for the Iris data set. For now, you just load all training data, fit the kNN model and predict the labels for a subset of the training data. Finally, you compare the results of the kNN model to the actual labels associated to the inputs.