

Series 5, May 3rd, 2016
(Probabilistic Modeling & Autoencoders)

It is not mandatory to submit solutions and sample solutions will be published in two weeks. If you choose to submit your solution, please send an e-mail from your `ethz.ch` address with subject `Exercise5` containing a PDF (\LaTeX or scan) to `lis2015@lists.inf.ethz.ch` until Monday, May 16th 2016.

Problem 1 (Independence Assumptions of Naive Bayes Classifiers):

Consider a naive Bayes classifier with binary class variable $C \in \{0, 1\}$ and two binary features $X_1 \in \{0, 1\}$ and $X_2 \in \{0, 1\}$. Assume that X_1 and X_2 are truly independent. You are given the following probabilities:

$$\begin{aligned}P(X_1 = 1|C = 1) &= p \\P(X_1 = 1|C = 0) &= 1 - p \\P(X_2 = 0|C = 1) &= q \\P(X_2 = 0|C = 0) &= 1 - q \\P(C = 0) &= P(C = 1) = 0.5\end{aligned}$$

- (a) Given a test sample with $X_1 = 1$ and $X_2 = 0$, compute the decision rule for classifying the example as belonging to class 1 in terms of q and p . Reformulate the decision rule in the form $p \geq \dots$.
- (b) We extend the naive Bayes classifier by adding another feature X_3 which is simply a copy of X_2 . Again, compute the decision rule of the classifier in terms of q and p . Reformulate the decision rule in the form $p \geq \dots$.
- (c) Compare the decision boundaries of (a) and (b) by varying the value of q between 0 and 1. Show where the second rule makes mistakes relative to the first (correct) decision rule.

Solution 1 (Independence Assumptions of Naive Bayes Classifiers):

- (a) Exploiting the independence assumption we compute

$$P(X_1 = 1, X_2 = 0, C = 1) = P(X_1 = 1, X_2 = 0|C = 1)P(C = 1) \tag{1}$$

$$= P(X_1 = 1|C = 1)P(X_2 = 0|C = 1)P(C = 1) \tag{2}$$

$$= pq\frac{1}{2}, \tag{3}$$

and, similarly,

$$P(X_1 = 1, X_2 = 0, C = 0) = P(X_1 = 1|C = 0)P(X_2 = 0|C = 0)P(C = 0) \tag{4}$$

$$= (1 - p)(1 - q)\frac{1}{2}. \tag{5}$$

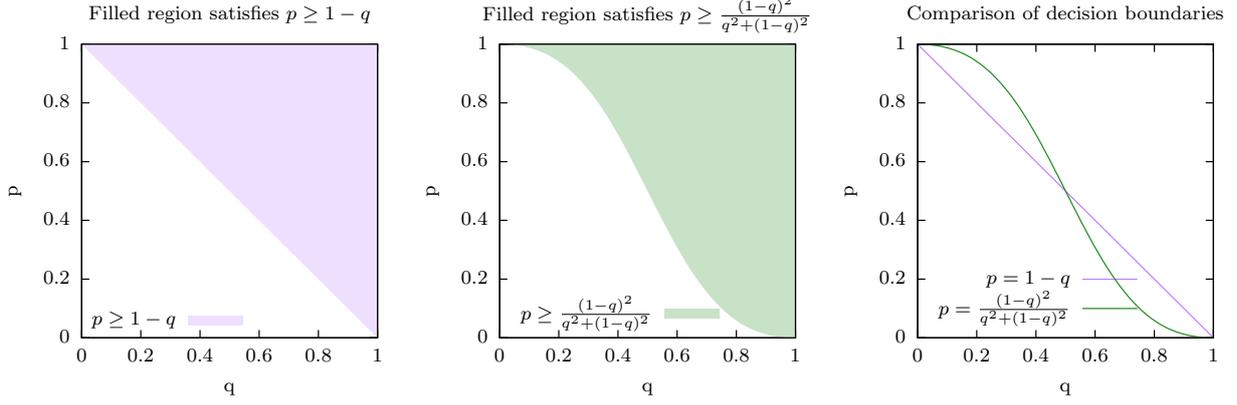


Figure 1: Decision boundaries

Consequently, we want

$$pq \frac{1}{2} \geq (1-p)(1-q) \frac{1}{2}. \quad (6)$$

Solving for p yields

$$p \geq 1 - q. \quad (7)$$

(b) Similarly, using the replicated feature one obtains

$$pqq \frac{1}{2} \geq (1-p)(1-q)(1-q) \frac{1}{2} \quad (8)$$

Solving for p yields

$$p \geq \frac{(1-q)^2}{q^2 + (1-q)^2}. \quad (9)$$

(c) The two boundaries are shown in Figure 1. Furthermore, the figure shows where one decision rule differs relative to the other (this is the region between the purple and the green line in the rightmost figure).

Problem 2 (Bayesian optimal decisions for logistic regression):

Apply Bayesian decision theory to derive the optimal decision rule for logistic regression in the following setting:

- Estimated conditional distribution: $\hat{P}(y|\mathbf{x}) = \begin{cases} \sigma(\mathbf{w}^T \mathbf{x}) & \text{if } y = 1 \\ 1 - \sigma(\mathbf{w}^T \mathbf{x}) & \text{if } y = -1 \end{cases}$
- Action set: $\{+1, -1, D\}$
- Cost function: $C(y, a) = \begin{cases} \mathbf{1}[y \neq a] & \text{if } a \in \{+1, -1\} \\ c < 0.5 & \text{if } a = D \end{cases}$

Here, $\mathbf{1}[\cdot]$ denotes the indicator function.

Solution 2 (Bayesian optimal decisions for logistic regression):

In the Bayesian optimal decision framework we want to pick the action which minimizes the expected loss:

$$a^* = \operatorname{argmin}_{a \in \mathcal{A}} \mathbb{E}_y[C(y, a) | x] = \operatorname{argmin}_{a \in \mathcal{A}} \int C(y, a)p(y | x)dy$$

In our case

$$a^* = \operatorname{argmin}_{a \in \mathcal{A}} \sum_y C(y, a)p(y | x) = \operatorname{argmin}_{a \in \mathcal{A}} \sum_y (\mathbf{1}[y \neq a]\mathbf{1}[D \neq a] + c\mathbf{1}[D = a])p(y | x)$$

Let's analyze the expected cost of each action separately:

$$\mathbb{E}[C(y, a) | x] = \begin{cases} p(y = -1 | x) & \text{if } a = 1 \\ p(y = 1 | x) & \text{if } a = -1 \\ c & \text{if } a = D \end{cases}$$

We can conclude that we should only pick an action $a \in \{-1, 1\}$ if the probability of action $-a$ being correct is lower than than c !

$$a^* = \begin{cases} y & \hat{P}(-y | x) < c \\ D & \text{otherwise} \end{cases}$$

Problem 3 (Bayesian optimal decisions for regression with asymmetric costs):

Apply Bayesian decision theory to derive the optimal decision rule for linear regression in the following setting:

- Estimated conditional distribution: $\hat{P}(y|\mathbf{x}) = \mathcal{N}(y; \mathbf{w}^T \mathbf{x}, \sigma^2)$
- Action set: \mathbb{R}
- Cost function: $C(y, a) = c_1 \max(y - a, 0) + c_2 \max(a - y, 0)$

Here, c_1 and c_2 denote positive real valued constants.

Solution 3 (Bayesian optimal decisions for regression with asymmetric costs):

We aim to minimize the expected risk $\mathbb{E}[C(y, a)] = \int C(y, a) \hat{P}(y|\mathbf{x}) dy$. Differentiating with respect to a yields

$$\frac{\partial}{\partial a} \mathbb{E}[C(y, a)] = \frac{\partial}{\partial a} \int_{-\infty}^{\infty} C(y, a) \hat{P}(y|\mathbf{x}) dy \quad (10)$$

$$= \frac{\partial}{\partial a} \int_{-\infty}^{\infty} c_1 \max(y - a, 0) \hat{P}(y|\mathbf{x}) dy + \int_{-\infty}^{\infty} c_2 \max(a - y, 0) \hat{P}(y|\mathbf{x}) dy \quad (11)$$

$$= -c_1 \int_a^{\infty} \hat{P}(y|\mathbf{x}) dy + c_2 \int_{-\infty}^a \hat{P}(y|\mathbf{x}) dy \quad (12)$$

$$= -c_1 [1 - \Phi(a; \mathbf{w}^T \mathbf{x}, \sigma^2)] + c_2 \Phi(a; \mathbf{w}^T \mathbf{x}, \sigma^2), \quad (13)$$

where $\Phi(u; v, w)$ is the CDF of the normal distribution with mean v and variance w . Setting the derivative to zero and applying some algebraic manipulations yields

$$\Phi(a; \mathbf{w}^T \mathbf{x}, \sigma^2) = \frac{c_1}{c_1 + c_2}. \quad (14)$$

Noting that $\Phi(u; v, w) = \Phi((u - v)/\sqrt{w}; 0, 1)$, we can rewrite the above equation as

$$\Phi\left(\frac{a - \mathbf{w}^T \mathbf{x}}{\sigma}; 0, 1\right) = \frac{c_1}{c_1 + c_2}. \quad (15)$$

Applying the inverse CDF of the standard normal distribution Φ^{-1} , and rewriting the result in terms of a , yields the optimal action

$$a^* = \mathbf{w}^T \mathbf{x} + \sigma \Phi^{-1}\left(\frac{c_1}{c_1 + c_2}\right). \quad (16)$$

Thus, if for example the cost of underestimation c_1 is smaller than the cost of overestimation c_2 , i.e. $c_1 < c_2$, then $\Phi^{-1}\left(\frac{c_1}{c_1 + c_2}\right) < 0$. In this way the risk of overestimation is reduced.

Problem 4 (Optional) (Autoencoders and PCA):

In this exercise, we analyze dimensionality reduction using autoencoders with linear activation functions and relate them to principal component analysis (PCA). We consider the following setup: let $D = \{x_1, \dots, x_N\}$ be given inputs, with $x_i \in \mathbb{R}^n$. Let $X = [x_1, \dots, x_N] \in \mathbb{R}^{n \times N}$ be the matrix formed from the inputs. Assume that we compute p hidden activations for every input x_i example according to $h_i = \phi_1(W_1 x_i + b_1)$, where $\phi_1(\cdot)$ is an activation function applied element-wise, $W_1 \in \mathbb{R}^{p \times n}$ are the input weights, and $b_1 \in \mathbb{R}^p$ are biases. Note that we can express the computation of all hidden activations as $H = \phi_1(W_1 X + b_1 u^T)$, where u is a vector containing only ones of size N . For this analysis, assume that $\phi_1(x) = x$. Given the hidden activations H and output weights $W_2 \in \mathbb{R}^{n \times p}$ as well as biases $b_2 \in \mathbb{R}^n$, the output of the autoencoder is computed as $Y = \phi_2(W_2 H + b_2 u^T)$, where again we assume $\phi_2(x) = x$. The weights and biases of the autoencoder are selected as

$$\arg \min_{W_1, W_2, b_1, b_2} \|X - Y\|^2. \quad (17)$$

- (a) Consider the squared-error criterion given the hidden activations, i.e. $\|X - (W_2 H + b_2 u^T)\|^2$. Derive an expression for the biases b_2 in terms of X , H and W_2 . Substitute your expression into the error and rewrite it in the form $\|X' - W_2 H'\|^2$, where X' (H') depends only on X (H) and constants.
- (b) Compare the problem of minimizing $\|X' - W_2 H'\|^2$ with the problem of computing the PCA from the lecture. Read off the optimal W_2 and H' . They should be expressed up to an arbitrary non-singular linear transform given by a $p \times p$ matrix T .
- (c) Show that the obtained solution for H' can actually be generated by proper choices of W_1 and b_1 .
- (d) Comment on the relation of W_1 to W_2 .
- (e) Comment on the transformation of the input computed by the autoencoder with respect to PCA.

Solution 4 (Autoencoders and PCA):

This solution is closely based on *H. Bourlard, and Y. Kamp, Auto-association by multilayer perceptrons and singular value decomposition, Biological Cybernetics, Springer-Verlag, vol. 59, pp. 291–294 1988.*

- (a) Note that $\|A\|^2 = \text{tr}(AA^T)$, where $\text{tr}(B)$ denotes the trace of matrix B . Furthermore, note that $\text{tr}(\cdot)$ is a linear operator and that $\text{tr}(A) = \text{tr}(A^T)$. Hence,

$$\|X - (W_2 H + b_2 u^T)\|^2 = \text{tr}((X - (W_2 H + b_2 u^T))(X - (W_2 H + b_2 u^T))^T) \quad (18)$$

$$= \text{tr}(X X^T) - 2\text{tr}(X(W_2 H + b_2 u^T)^T) + \text{tr}((W_2 H + b_2 u^T)(W_2 H + b_2 u^T)^T) \quad (19)$$

$$= \text{tr}(X X^T) - 2\text{tr}(X H^T W_2) - 2\text{tr}(X u b_2^T) + \text{tr}(W_2 H H^T W_2) \quad (20)$$

$$+ 2\text{tr}(W_2 H u b_2^T) + \text{tr}(b_2 u^T u b_2^T) \quad (21)$$

Dropping all terms that do not depend on b_2 (these will cancel once we compute the gradient with respect to b_2) and noting that $u^T u = N$, we get

$$-2\text{tr}(X u b_2^T) + 2\text{tr}(W_2 H u b_2^T) + N\text{tr}(b_2 b_2^T). \quad (22)$$

Computing the gradient of the above expression with respect to b_2 yields

$$\nabla_{b_2} [-2\text{tr}(X u b_2^T) + 2\text{tr}(W_2 H u b_2^T) + N\text{tr}(b_2 b_2^T)] = -2\nabla_{b_2} \text{tr}(X u b_2^T) + 2\nabla_{b_2} \text{tr}(W_2 H u b_2^T) + N\nabla_{b_2} \text{tr}(b_2 b_2^T) \quad (23)$$

$$= -2X u + 2W_2 H u + 2N b_2. \quad (24)$$

Equating to zero and solving for b_2 yields

$$b_2 = \frac{1}{N}(X - W_2H)u. \quad (25)$$

Substituting this result into $\|X - (W_2H + b_2u^T)\|^2$ gives

$$\|X - (W_2H + b_2u^T)\|^2 = \|X - (W_2H + \frac{1}{N}(X - W_2H)uu^T)\|^2 \quad (26)$$

$$= \|X(I - \frac{1}{N}uu^T) - W_2[H(I - \frac{1}{N}uu^T)]\|^2 \quad (27)$$

$$= \|X' - W_2H'\|^2, \quad (28)$$

where $X' = X(I - \frac{1}{N}uu^T)$ and $H' = H(I - \frac{1}{N}uu^T)$.

- (b) Assume that $p < n$, i.e. W_2 has rank p . From the lecture we know that the optimal solution to $\|X' - W_2H'\|^2$ is given when the product W_2H' is the best rank p approximation of X' . This approximation can be obtained by computing the singular value decomposition (SVD) of X' , i.e.

$$X' = U\Sigma V^T, \quad (29)$$

where U and V are matrices of size $n \times n$ and $N \times n$, respectively, consisting of the normalized eigenvectors of $X'X'^T$ and $X'^T X'$, respectively. Σ is a diagonal matrix with entries $\sqrt{\lambda_1}, \dots, \sqrt{\lambda_n}$ on the diagonal, where λ_i is the i -th largest eigenvalue (assume the eigenvalues to be sorted such that $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$). In this setting, the best rank p approximation of X' is given by

$$W_2H' = \tilde{U}\tilde{\Sigma}\tilde{V}^T, \quad (30)$$

where $\tilde{\Sigma}$ is a diagonal matrix with the p largest eigenvalues $\lambda_1, \dots, \lambda_p$ as its entries, and where \tilde{U} and \tilde{V}^T are formed by the first p columns of U and V , respectively.

Consequently, denoting by T an arbitrary invertible $p \times p$ matrix, we have $W_2 = \tilde{U}T^{-1}$ and $H' = T\tilde{\Sigma}\tilde{V}^T$.

- (c) It remains to show that the H' can actually be generated by proper choices of W_1 and b_1 . The output H of the first layer is computed according to $H = W_1X + b_1u^T$. Multiplying from the right both sides by $(I - \frac{1}{N}uu^T)$ gives

$$H(I - \frac{1}{N}uu^T) = H' = W_1X' + b_1u^T(I - \frac{1}{N}uu^T). \quad (31)$$

Hence, we need that

$$T\tilde{\Sigma}\tilde{V}^T = H' \quad (32)$$

$$= W_1X' + b_1u^T(I - \frac{1}{N}uu^T). \quad (33)$$

Observe that

$$b_1u^T(I - \frac{1}{N}uu^T) = b_1u^T - \frac{1}{N}b_1u^Tuu^T \quad (34)$$

$$= b_1u^T - \frac{1}{N}b_1Nu^T \quad (35)$$

$$= b_1u^T - b_1u^T \quad (36)$$

$$= 0, \quad (37)$$

and hence w_1 is arbitrary. Furthermore,

$$W_1X' = T\tilde{\Sigma}\tilde{V}^T \quad (38)$$

$$= T\tilde{U}^T X', \quad (39)$$

where we used $X' = U\Sigma V^T$. Consequently, $W_1 = T\tilde{U}^T$ is a proper choice for W_1 .

- (d) We found that $W_1 = T\tilde{U}^T$ and $W_2 = \tilde{U}T^{-1}$. Setting $T = I$, we have $W_2 = W_1^T$.
- (e) An auto-encoder essentially computes the PCA (up to an invertible linear transform).