

1 Kernels

Here we will look at an example of constructing feature spaces for classification problems. We will explore the idea of the "kernel trick" applied to classifiers other than SVM (i.e., classifiers using loss functions other than hinge loss). Specifically, we will introduce the kernel logistic regression (KLR). Recall that in class we have described the general logistic regression model as

$$\hat{p}(y = 1 | \mathbf{x}; \mathbf{w}, w_0) = \frac{1}{1 + \exp\left(-\sum_{j=1}^d w_j \phi_j(\mathbf{x})\right)},$$

where $\phi_j(\mathbf{x})$ is the j -th basis function, or feature - generally, a function mapping $\mathcal{X} \rightarrow \mathbb{R}$.

Problem 1 [15 points]

Show how by appropriate choice of basis functions ϕ , given a training set $\mathbf{x}_1, \dots, \mathbf{x}_N$, one can obtain a logistic regression model whose predictions on a test point \mathbf{x}_0 depend on the training data only through the kernel values $K(\mathbf{x}_i, \mathbf{x}_0)$ for $i = 1, \dots, N$. Then write down the gradient of the loss, with L_2 regularization, on a single example, and show that the training for this model via gradient descent also depends on the training data only through kernel computations.

End of problem 1

Let $\phi_j(\mathbf{x}) = \text{ker}(\mathbf{x}_j, \mathbf{x})$, $j \in [1, N]$. Let $d=N$.

Then

$$\hat{p}(y=1 | \mathbf{x}; \mathbf{w}, w_0) = \frac{1}{1 + \exp\left(-\sum_{j=1}^d w_j k(\mathbf{x}_j, \mathbf{x}_0)\right)}$$

$$\text{Let } \theta = \sum_{j=1}^d w_j k(\mathbf{x}_j, \mathbf{x}_0).$$

For one sample,

$$L(\mathbf{w}) = -\left(y_i \cdot \log \frac{e^\theta}{1+e^\theta} + (1-y_i) \log \left(1 - \frac{e^\theta}{1+e^\theta}\right)\right) + \sum_{j=1}^d \lambda \|\mathbf{w}_j\|^2$$

$$= -\left(-y_i \cdot \log(1+e^\theta) + \theta - \log(1+e^\theta) - y_i \theta + y_i \log(1+e^\theta)\right) + \lambda \sum_{j=1}^d \|\mathbf{w}_j\|^2$$

$$= -\left(\theta - y_i \theta - \log(1+e^\theta)\right) + \lambda \sum_{j=1}^d \|\mathbf{w}_j\|^2$$

$$= -\left(-y_i \theta + \log e^\theta - \log(1+e^\theta)\right) + \lambda \sum_{j=1}^d \|\mathbf{w}_j\|^2$$

$$= -\left(y_i \theta + \log\left(\frac{e^\theta}{1+e^\theta}\right)\right) + \lambda \sum_{j=1}^d \|\mathbf{w}_j\|^2$$

$$\log \frac{e^\theta}{1+e^\theta} = \log \frac{1}{e^{-\theta} + 1}$$

$$= -(y_i \theta - \log(1+e^{-\theta})) + \lambda \sum_{j=1}^d \|\mathbf{w}_j\|^2$$

Take derivative,

$$\frac{\partial}{\partial \mathbf{w}_j} L(\mathbf{w}) = -\left(y_i k(\mathbf{x}_j, \mathbf{x}_0) - \frac{e^{-\theta} \cdot -k(\mathbf{x}_j, \mathbf{x}_0)}{1+e^{-\theta}}\right)$$

$$+ 2\lambda \|\mathbf{w}_j\|$$

$$= -y_i k(\mathbf{x}_j, \mathbf{x}_0) - \frac{k(\mathbf{x}_j, \mathbf{x}_0)}{1+e^\theta} + 2\lambda \|\mathbf{w}_j\|$$

Therefore, the gradient is

$$\nabla L(\mathbf{w}) = \begin{pmatrix} \frac{\partial L}{\partial \mathbf{w}_1} \\ \vdots \\ \frac{\partial L}{\partial \mathbf{w}_d} \end{pmatrix}$$

with the update rule

$$\mathbf{w}_{\text{new}} = \mathbf{w} - \alpha \cdot \nabla L$$

We observe that the process depends only on training data through kernel.

$$\log \frac{1}{1+e^\theta}$$

2 Support Vector Machines

Now we will consider some details of the dual formulation of SVM, the one in which we optimize the Lagrange multipliers α_i . In class we saw how to derive a constrained quadratic program, which can then be "fed" to an off-the-shelf quadratic program solver. These solvers are usually constructed to handle certain standard formulations of the objective and constraints.

Specifically the canonical form of a quadratic program with linear constraints is, mathematically:

$$\operatorname{argmin}_{\alpha} \frac{1}{2} \alpha^T H \alpha + f^T \alpha, \quad (1)$$

$$\text{such that: } A \cdot \alpha \leq a, \quad (2)$$

$$B \cdot \alpha = b, \quad (3)$$

$$(4)$$

The vector $\alpha \in \mathbb{R}^N$, where N is the number of training examples, contains the unknown variables to be solved for. The matrix $H \in \mathbb{R}^{N \times N}$ and vector $f \in \mathbb{R}^N$ specify the quadratic objective; the matrix $A \in \mathbb{R}^{k_{\text{ineq}} \times N}$ and vector $a \in \mathbb{R}^{k_{\text{ineq}}}$ specify k_{ineq} inequality constraints. Similarly, $B \in \mathbb{R}^{k_{\text{eq}} \times N}$ and vector $b \in \mathbb{R}^{k_{\text{eq}}}$ specify k_{eq} equality constraints. Note that you can express a variety of inequality constraints by adding rows to B and elements to b ; think how you would do it to express, e.g., a "greater or equal" constraint.

Problem 2 [15 bonus points]

Describe in detail how you would compute H , f , A , a , B , and b , to set up the dual optimization problem for the kernel SVM

$$\operatorname{argmin}_w \left\{ \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N \max[0, 1 - y_i (w^T \phi(x_i) - w_0)] \right\}$$

given a kernel function $K(\cdot, \cdot)$ corresponding to the dot product in ϕ space, and N training examples (x_i, y_i)

End of problem 2

Approach 1:

• use z_i to represent $\max\{0, 1 - y_i (w^T \phi(x_i) - w_0)\}$,
with $\begin{cases} z_i \geq 0 \\ z_i \geq 1 - y_i (w^T \phi(x_i) - w_0) \end{cases}$.

• assign $\begin{cases} \alpha_i \cdot z_i \leq 0, \\ \beta_i \cdot z_i \leq 1 - y_i (w^T \phi(x_i) - w_0) \end{cases}$, where $\alpha_i \in \mathbb{U}$, $\beta_i \leq 0$.

• Then, let

$$H = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}, \quad d = \begin{pmatrix} w \\ z \end{pmatrix}, \quad f = \begin{pmatrix} 0 \\ 0 \\ c \\ 1 \end{pmatrix} \begin{pmatrix} w \\ z \end{pmatrix}$$

We get

$$\operatorname{argmax}_{\alpha_i, \beta_i} \left\{ \frac{1}{2} \cdot (w, z) \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} w \\ z \end{pmatrix} + \underbrace{(0, \dots, 0, c, \dots, c)}_d \cdot \underbrace{\begin{pmatrix} w \\ z \end{pmatrix}}_N \right\}$$

$$= \operatorname{argmax}_{\alpha_i, \beta_i} \left\{ \frac{1}{2} \|w\|^2 + C \cdot \sum_{i=1}^N z_i \right\}$$

Approach 2:

• From Tutorial 5, the dual is

$$\max_{d_i} \sum_{i=1}^N d_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N d_i d_j y_i y_j K(x_i) \cdot \phi(x_j)$$

$$\begin{cases} 0 \leq d_i \leq C \\ \sum_{i=1}^N d_i y_i = 0 \end{cases}$$

which is

$$\min_d \left\{ \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N d_i d_j y_i y_j K(x_i) \cdot \phi(x_j) - \sum_{i=1}^N d_i \right\}$$

$$\begin{cases} 0 \leq d_i \leq C \\ \sum_{i=1}^N d_i y_i = 0 \end{cases}$$

• Define

$$f = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \Bigg\}_N$$

$$a = \begin{pmatrix} C \\ 0 \\ \vdots \\ 0 \end{pmatrix} \Bigg\}_N$$

$$A = \begin{pmatrix} I_N \\ -I_N \end{pmatrix} \Bigg\}_{N \times N}$$

$$b = 0$$

$$B = (y_1, \dots, y_N)$$

$$H = \begin{pmatrix} 0 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{pmatrix} \Bigg\}_N$$

$H_{ij} = y_i y_j K(x_i, x_j)$

This will reconstruct our dual function.

Next, we will consider finding the value of b , the threshold (bias) term for a kernel SVM

$$h(\mathbf{x}) = \text{sign} \left(\sum_{i: \alpha_i > 0} \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b \right),$$

where α_i are the coefficients for the support vectors in the training data.

Problem 3 [10 points]

Suppose you have solved for α_i in the SVM classifier. Explain how exactly you can calculate b .

End of problem 3

• Let C be the margin length.

$$w = \sum_{i=1}^N d_i y_i \phi(x_i)$$



As only the points exactly

on the boundary matter,

we consider the points with $0 < d_i < C$.

• Pick x_j on the boundary,

$$y_i \left(\sum_{\substack{i \neq j \\ 0 < d_i < C}} d_i y_i K(x_i, x_j) + b \right) = 1$$

$$b = \frac{1}{y_j} - \sum_{\substack{i \neq j \\ 0 < d_i < C}} d_i y_i K(x_i, x_j)$$

$y_j = \pm 1$

4 Multivariate Gaussians

Since many (probably most) generative models for real-valued data involve multivariate Gaussian distribution, it is worth getting to know them better. In particular, here we will understand why one always sees those elliptical contours drawn when Gaussians are visualized in 2D.

Recall that the probability density function (pdf) of a Gaussian distribution in \mathbb{R}^d is given by

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)\right). \quad (5)$$

Problem 6 [10 points]

Show that a contour corresponding to a fixed value of pdf is an ellipse in the 2D space, $\mathbf{x} = [x_1, x_2]$.

End of problem 6

Advice: You may find it easier to work in the log domain, and to write out explicitly the expression in (5) in terms of x_1 and x_2 .

$$\frac{1}{(2\pi)^d |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)\right) = \text{const}$$

$$(\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu) = \text{const}$$

$$\bullet \text{ case } \Sigma = \begin{pmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{pmatrix}, \quad \Sigma^{-1} = \begin{pmatrix} \frac{1}{\sigma_1^2} & 0 \\ 0 & \frac{1}{\sigma_2^2} \end{pmatrix}$$

$$(x_1 - \mu_1, x_2 - \mu_2) \begin{pmatrix} \frac{1}{\sigma_1^2} & 0 \\ 0 & \frac{1}{\sigma_2^2} \end{pmatrix} \begin{pmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{pmatrix} = C$$

$$\frac{(x_1 - \mu_1)^2}{\sigma_1^2} + \frac{(x_2 - \mu_2)^2}{\sigma_2^2} = C$$

Divide C on both sides, we get the form of an ellipse.

$$\bullet \text{ case } \Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{12}^2 \\ \sigma_{12}^2 & \sigma_2^2 \end{pmatrix},$$

As covariance matrices are symmetric, they are diagonalizable.

write

$$\Sigma = P D P^{-1}$$

then

$$\Sigma^{-1} = P D^{-1} P^{-1}$$

$$(x_1 - \mu_1, x_2 - \mu_2) \begin{pmatrix} a & b \\ c & d \end{pmatrix} \begin{pmatrix} u_1 & 0 \\ 0 & u_2 \end{pmatrix} \begin{pmatrix} a' & b' \\ c' & d' \end{pmatrix} \begin{pmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{pmatrix}$$

$$\stackrel{\textcircled{1}}{=} \underbrace{(x_1 - \mu_1) \begin{pmatrix} a & b \\ c & d \end{pmatrix} \begin{pmatrix} u_1 & 0 \\ 0 & u_2 \end{pmatrix}}_{\textcircled{1}} \underbrace{\begin{pmatrix} a' & b' \\ c' & d' \end{pmatrix} \begin{pmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{pmatrix}}_{\textcircled{2}}$$

$$= [(x_1 - \mu_1) a \cdot u_1, 0] + [0, (x_2 - \mu_2) d u_2]$$

$$= [(x_1 - \mu_1) a u_1, (x_2 - \mu_2) d u_2]$$

$$\stackrel{\textcircled{2}}{=} \begin{bmatrix} a'(x_1 - \mu_1) + b'(x_2 - \mu_2) \\ c'(x_1 - \mu_1) + d'(x_2 - \mu_2) \end{bmatrix}$$

multiply $\textcircled{1}$ and $\textcircled{2}$:

We should obtain the form

$$\frac{(x_1 - \mu_1)^2}{c_1^2} + \frac{(x_2 - \mu_2)^2}{c_2^2} = 1$$

5 Generative models

With help from Haichuan Wang

Here we will consider the Gaussian generative model for $\mathbf{x} \in \mathbb{R}^d$ which assumes equal, isotropic covariance matrices for each class:

$$\Sigma_c = \begin{bmatrix} \sigma_1^2 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma_d^2 \end{bmatrix}. \quad (6)$$

The classes of course have separate means. That is, the conditional density of the j -th coordinate of \mathbf{x} given class c is a Gaussian $N(\mu_{c,j}, \sigma_j^2)$, with these densities independent (given class) for different values of j .

For simplicity, let's consider a two-class problem, with $y \in \{0, 1\}$. As we discussed in class, when training the generative model, we simply fit $p(\mathbf{x}|y)$ and $p(y)$ to the training data, and use the resulting discriminant analysis to produce a decision rule which predicts

$$\hat{y}(\mathbf{x}) = \operatorname{argmax}_c \{p(\mathbf{x}|y=c)p(y=c)\}. \quad (7)$$

However, this model does also produce an (implicit) estimate for the posterior $p(y=c|\mathbf{x})$.

Problem 7 [15 points]

Show that the posterior $p(y=c|\mathbf{x})$ resulting from the generative model above has the same form as the posterior in logistic regression model,

$$p(y=c|\mathbf{x}) = \frac{1}{1 + \exp(b + \mathbf{w} \cdot \mathbf{x})}. \quad (8)$$

for appropriate values of $b \in \mathbb{R}$, $\mathbf{w} \in \mathbb{R}^d$.

End of problem 7

Advice: Start with Bayes rule, and use the simplifying assumptions in (6) to derive the posterior.

• Bayes rule:

$$p(y=c|y) = \frac{p(\mathbf{x}|y=c)p(y=c)}{p(\mathbf{x})} = \frac{p(\mathbf{x}|y=c) \cdot p(y=c)}{\dots}$$

• As each class has $N(\mu_{c,j}, \sigma_j^2)$ distribution,

$$\begin{aligned} p(\mathbf{x}|y=c) &= \frac{1}{(2\pi)^{\frac{d}{2}} (\sigma_1^2 \cdots \sigma_d^2)^{\frac{1}{2}}} \cdot \exp\left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_c)^T \begin{pmatrix} \frac{1}{\sigma_1^2} & & \\ & \ddots & \\ & & \frac{1}{\sigma_d^2} \end{pmatrix} (\mathbf{x} - \boldsymbol{\mu}_c)\right) \\ &= \prod_{j=1}^d \frac{1}{\sqrt{2\pi} \sigma_j} \exp\left(-\frac{1}{2\sigma_j^2} (x_j - \mu_{c,j})^2\right) \end{aligned}$$

$$p(\mathbf{x}|y=c') = \prod_{j=1}^d \frac{1}{\sqrt{2\pi} \sigma_j} \exp\left(-\frac{1}{2\sigma_j^2} (x_j - \mu_{c',j})^2\right)$$

• Now,

$$p(y=c|\mathbf{x}) = \frac{p(\mathbf{x}|y=c) \cdot p(y=c)}{p(\mathbf{x}|y=c) \cdot p(y=c) + p(\mathbf{x}|y=c') \cdot p(y=c')}$$

$$= \frac{1}{1 + \frac{p(\mathbf{x}|y=c') p(y=c')}{p(\mathbf{x}|y=c) p(y=c)}}$$

$$= \frac{1}{1 + \frac{p(y=c')}{p(y=c)} \cdot \exp\left(\sum_{j=1}^d \frac{1}{2\sigma_j^2} ((x_j - \mu_{c',j})^2 - (x_j - \mu_{c,j})^2)\right)}$$

$$= \frac{1}{1 + \exp\left(\sum_{j=1}^d \frac{\mu_{c,j} - \mu_{c',j}}{\sigma_j^2} x_j - \sum_{j=1}^d \frac{1}{2\sigma_j^2} (\mu_{c,j}^2 + \mu_{c',j}^2) (\mu_{c,j} - \mu_{c',j})\right) + (\log p(y=c') - \log p(y=c))}$$

• To get the form of

$$\frac{1}{1 + \exp(b + \mathbf{w} \cdot \mathbf{x})},$$

We can let

$$\mathbf{w} = \begin{pmatrix} \frac{\mu_{c,1} - \mu_{c',1}}{\sigma_1^2} \\ \vdots \\ \frac{\mu_{c,d} - \mu_{c',d}}{\sigma_d^2} \end{pmatrix}$$

$$b = \log p(y=c') - \log p(y=c) - \sum_{j=1}^d \frac{1}{2\sigma_j^2} (\mu_{c,j}^2 + \mu_{c',j}^2) (\mu_{c,j} - \mu_{c',j})$$

Problem 8 [10 points]

The previous problem established that the two models – logistic regression and the linear discriminant analysis based on the isotropic Gaussian model (6) – have the same form of the posterior $p(y|x)$. Will the two models produce the same classifier when applied to a given training set? Why or why not?

End of problem 8

- ' The two models will not produce the same classifier, because in the Gaussian model, we made many additional assumptions .e.g. the distribution being $N(\mu_c, \Sigma_c)$ for each class, the coordinates are independent, etc.
- ' However, for the logistic regression, the assumptions are less. Therefore, the generative model is harder to fit a distribution when there are more data. But logistic regression performs better when the training set is large.