

# CPSC 340 Assignment 5 (due Friday November 15 at 11:55pm)

If you are planning to do a poster project, make sure to register for the poster session as soon as possible. Details regarding registration are here:  
<https://piazza.com/class/k02x04b6o524?cid=358>

## 1 Kernel Trick

In this question you will revisit questions from previous assignments, this time implementing the same (or similar) models using the “kernel trick”.

### 1.1 “Other” Normal Equations

The script *example\_nonLinear* loads a dataset from a previous assignment, and fits an L2-regularized least squares model with a bias term (the regularization leads to small improvement in the test error, even for this 2-variable problem). Modify the *leastSquaresBiasL2* function so that it uses the “other” normal equations we discussed in class,

$$v = Z^T \underbrace{(ZZ^T + \lambda I)^{-1} y}_u.$$

In particular, the training should result in an  $n \times 1$  vector  $u$  of parameters, and predictions are made based on the training  $Z$  and the vector  $u$ . [Hand in your code for the modified function.](#)

Hint: you should get the same predictions as the original function. To help debugging, you may want to first re-write the calculation of  $v$  using the above formula, to verify that this gives you the same vector  $v$ .

```
function leastSquaresBiasL2(X,y,lambda)

    %# Add bias column
    n = size(X,1)
    Z = [ones(n,1) X]

    %# Find regression weights minimizing squared error
    v1 = (Z'*Z + lambda*I)\(Z'*y)
    v = Z'*(Z'*Z + lambda * I)^-1*y
    @show(v1, v)
    %# Make linear prediction function
    predict(Xhat) = [ones(size(Xhat,1),1) Xhat]*v

    %# Return model
    return LinearModel(predict,v)
end
```

Answer:

### 1.2 Polynomial Kernel

Write a new function, *leastSquaresKernelBasis*, taking a degree  $p$  and a regularization parameter  $\lambda$ . It should fit a degree- $p$  polynomial to the data, using the kernel trick to avoid ever forming the matrix  $Z$ . [Hand in your code and the plot obtained with  \$p = 3\$  and  \$\lambda = 10^{-6}\$ .](#)

Hint: you may find it helpful to write a function *polyKernel* that takes two matrices as inputs (either  $X$  and  $X$  or  $\tilde{X}$  and  $X$ ) and a degree  $p$ , and computes the polynomial kernel between all pairs of rows in the matrices.

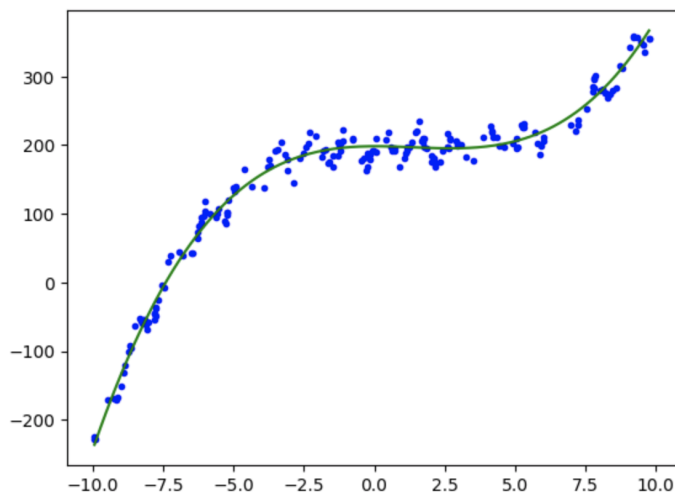
Answer:

```
function leastSquaresKernelBias(X,y,lambda, p)

—# Find regression weights minimizing squared error
—k = polyKernel(X,X, p)
—# Make linear prediction function
—predict(Xhat) = polyKernel(Xhat, X, p) * ((k + lambda * I) \ y)

—# Return model
—return GenericModel(predict)
end

function polyKernel(X1, X2, p)
    t = size(X1,1)
    n = size(X2,1)
    k = ones(t,n)
    for i = 1:t
        for j = 1:n
            k[i,j] = (1 + X1[i,:]'*X2[j,:])^p
        end
    end
    return k
end
```



Squared train Error with least squares: 252.017  
Squared test Error with least squares: 242.791

### 1.3 Gaussian-RBF Kernel

Repeat the previous question and report the same quantities, but using the Gaussian RBF kernel. You can use  $\sigma = 1$  and  $\lambda = 10^{-6}$  for the plot.

Answer:

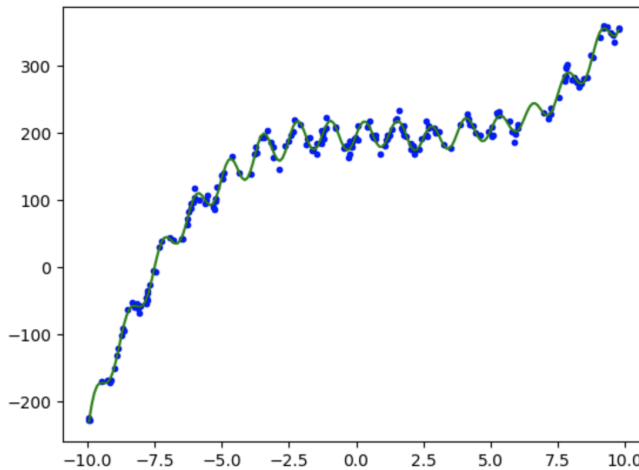
```

function RBFKernelBias(X,y,lambda, sigma)
    # Find regression weights minimizing squared error
    k = RBFKernel(X,X, sigma)
    # Make linear prediction function
    predict(Xhat) = RBFKernel(Xhat, X, sigma) * ((k + lambda * I) \ y)

    # Return model
    return GenericModel(predict)
end

function RBFKernel(X1, X2, sigma)
    t = size(X1,1)
    n = size(X2,1)
    k = ones(t,n)
    for i = 1:t
        for j = 1:n
            k[i,j] = exp(- norm(X1[i,:] - X2[j,:])^2/(2*sigma^2))
        end
    end
    return k
end

```



Squared train Error with least squares: 39.163  
 Squared test Error with least squares: 70.580

1]: 1-element Array{PyCall.PyObject,1}:  
 PyObject <matplotlib.lines.Line2D object at 0x14333eac8>

## 2 MAP Estimation

In class, we considered MAP estimation in a regression model where we assumed that:

- The likelihood  $p(y_i | x_i, w)$  for each example  $i$  is a normal distribution with a mean of  $w^T x_i$  and a variance of 1.
- The prior  $p(w_j)$  for each variable  $j$  is a normal distribution with a mean of zero and a variance of  $\lambda^{-1}$ .

Under these assumptions we showed that computing the MAP estimate with  $n$  training examples leads to the standard L2-regularized least squares objective function:

$$f(w) = \frac{1}{2} \|Xw - y\|^2 + \frac{\lambda}{2} \|w\|^2.$$

For each of the alternate assumptions below, write down the objective function that results (from minimizing the negative log-posterior, and simplifying as much as possible):

1. We use a Laplace likelihood with a mean of  $w^T x_i$  and a scale of 1, and we use a zero-mean Laplace

prior for each variable with a scale parameter of  $\lambda^{-1}$ ,

$$p(y_i | x_i, w) = \frac{1}{2} \exp(-|w^T x_i - y_i|), \quad p(w_j) = \frac{\lambda}{2} \exp(-\lambda |w_j|).$$

Answer:

$$\begin{aligned} f(w) &= - \sum_{i=1}^n \log(p(y_i | x_i, w)) \\ &= - \sum_{i=1}^n \log\left(\frac{1}{2} \exp(-|w^T x_i - y_i|)\right) \\ &= - \sum_{i=1}^n \log\left(\frac{1}{2}\right) - \sum_{i=1}^n \log(\exp(-|w^T x_i - y_i|)) \\ &= \text{constant} - \sum_{i=1}^n |w^T x_i - y_i| \\ &= \text{constant} + \|w^T x - y\|_1 \end{aligned}$$

2. We use a normal likelihood with a mean of  $w^T x_i$  but where each example  $i$  has its own positive variance  $\sigma_i^2$ , and a normal prior with a variance of  $\lambda^{-1}$  and a mean that is some “guess”  $w^0$  of the optimal parameter vector,

$$p(y_i | x_i, w) = \frac{1}{\sqrt{2\sigma_i^2\pi}} \exp\left(-\frac{(w^T x_i - y_i)^2}{2\sigma_i^2}\right), \quad p(w_j) \propto \exp\left(-\frac{\lambda(w_j - w_j^0)^2}{2}\right).$$

Answer:

$$\begin{aligned} f(w) &= - \sum_{i=1}^n \log(p(y_i | x_i, w)) \\ &= - \sum_{i=1}^n \log\left(\frac{1}{\sqrt{2\sigma_i^2\pi}} \exp\left(-\frac{(w^T x_i - y_i)^2}{2\sigma_i^2}\right)\right) \\ &= - \sum_{i=1}^n \left(\log\left(\frac{1}{\sqrt{2\sigma_i^2\pi}}\right) + \log\left(\exp\left(-\frac{(w^T x_i - y_i)^2}{2\sigma_i^2}\right)\right)\right) \\ &= - \sum_{i=1}^n \left(\text{constant} + \left(-\frac{(w^T x_i - y_i)^2}{2\sigma_i^2}\right)\right) \\ &= \text{constant} + \|\Sigma^{-1}(w^T x - y)\|_2^2 \end{aligned}$$

The standard notation for this case is to use  $\Sigma$  as a diagonal matrix with the  $\sigma_i^2$  values along the diagonal.

3. We use a Poisson likelihood with a mean of  $\exp(w^T x_i)$ ,<sup>1</sup> and we use a uniform prior for some constant  $\kappa$ ,

$$p(y_i | x_i, w) = \frac{\exp(y_i w^T x_i) \exp(-\exp(w^T x_i))}{y_i!}, \quad p(w_j) \propto \kappa$$

---

<sup>1</sup>This is one way to use regression to model *counts*, like “number of Facebook likes”.

Answer:

$$\begin{aligned}
 f(w) &= - \sum_{i=1}^n \log(p(y_i | x_i, w)) \\
 &= - \sum_{i=1}^n \log\left(\frac{\exp(y_i w^T x_i) \exp(-\exp(w^T x_i))}{y_i!}\right) \\
 &= - \sum_{i=1}^n y_i w^T x_i + -\exp(w^T x_i) - \log(y_i!)
 \end{aligned}$$

For this sub-question you don't need to put likelihood in matrix notation.

4. We use a Laplace likelihood with a mean of  $w^T x_i$  where each example  $i$  has its own positive scale parameter  $v_i^{-1}$ , and a student  $t$  prior (which is very robust to irrelevant features) with  $\nu$  degrees of freedom,

$$p(y_i | x_i, w) = \frac{1}{2} \exp(-v_i |w^T x_i - y_i|), \quad p(w_j) = \frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\nu\pi} \Gamma(\frac{\nu}{2})} \left(1 + \frac{w_j^2}{\nu}\right)^{-\frac{\nu+1}{2}}$$

where you use can  $V$  as a diagonal matrix with the  $v_i$  along the diagonal and  $\Gamma$  is the “gamma” function (which is always non-negative). You do not need to put the log-prior in matrix notation.

Answer:

$$\begin{aligned}
 f(w) &= - \sum_{i=1}^n \log(p(y_i | x_i, w)) \\
 &= - \sum_{i=1}^n \log\left(\frac{1}{2} \exp(-v_i |w^T x_i - y_i|)\right) \\
 &= - \sum_{i=1}^n \log(1/2) + \log(\exp(-v_i |w^T x_i - y_i|)) \\
 &= \text{constant} - \sum_{i=1}^n (-v_i |w^T x_i - y_i|) \\
 &= \text{constant} + \|(w^T X - y)V\|_1
 \end{aligned}$$

## 3 Principal Component Analysis

### 3.1 PCA by Hand

Consider the following dataset, containing 5 examples with 2 features each:

$x_1$	$x_2$
-4	3
0	1
-2	2
4	-1
2	0

Recall that with PCA we usually assume that the PCs are normalized ( $\|w\| = 1$ ), that we need to center the data before we apply PCA, and that the direction of the first PC is the one that minimizes the orthogonal distance to all data points.

1. What is the first principal component?

Answer: Normalize X:  $(X_2 - 1)$

$x_1$	$x_2$	
-4	2	
0	0	
-2	1	Then we can see the line is $x_2 = -\frac{x_1}{2}$ .
4	-2	
2	-1	

We then can have  $W = [-2, 1]^T$

2. What is the (L2-norm) reconstruction error of the point (-3, 2.5)? (Show your work.)

Answer: We first do the "normalize". We get data: (-3,1.5)

Thus We can see that reconstruction is  $1.5 * [-2, 1]^T + [0, 1]$

3. What is the (L2-norm) reconstruction error of the point (-3, 2)? (Show your work.)

Answer: We first do the "normalize". We get data: (-3,1) The slope of line that is certical to  $x_2 = -\frac{x_1}{2}$  is  $2x_1 = x_2$ . We plug (-3,1) in to the equation

$2x_1 = x_2 + c$  we can resolve that  $c = -8$ . We then would get the intersection of the current line to w. We have

$$\begin{cases} x_2 = -\frac{x_1}{2} \\ 2x_1 = x_2 - 7 \end{cases}$$

$$\begin{cases} x_1 = -\frac{14}{5} \\ x_2 = \frac{7}{5} \end{cases}$$

Thus the reconstruction is  $\frac{7}{5} * [-2, 1]^T + [0, 1] = [-2.8, 2.4]$

## 3.2 Data Visualization

The script *example\_PCA* will load a dataset containing 50 examples, each representing an animal. The 85 features are traits of these animals. The script standardizes these features and gives two unsatisfying visualizations of it. First it shows a plot of the matrix entries, which has too much information and thus gives little insight into the relationships between the animals. Next it shows a scatterplot based on two random features and displays the name of 10 randomly-chosen animals. Because of the binary features even a scatterplot matrix shows us almost nothing about the data.

The function *PCA* applies the classic PCA method (orthogonal bases via SVD) for a given  $k$ . Using this function, modify the demo so that the scatterplot uses the latent features  $z_i$  from the PCA model with  $k = 2$ . Make a scatterplot of the two columns in  $Z$ , and use the *annotate* function to label a bunch of the points in the scatterplot.

1. Hand in your modified demo and the scatterplot.

Answer:

```

using DelimitedFiles

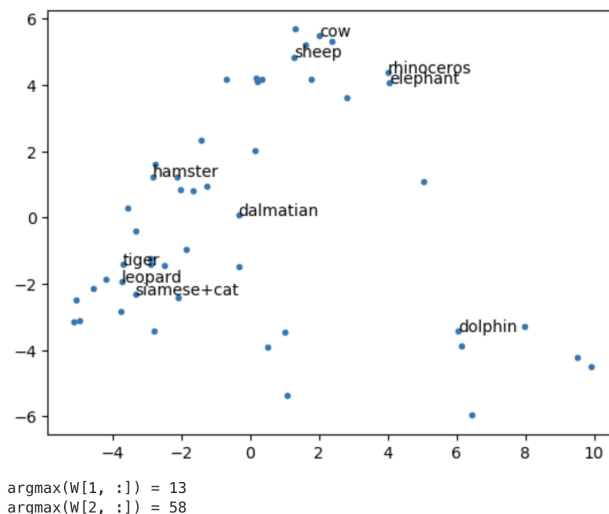
# Load data
dataTable = readdlm("animals.csv",',')
X = float(real(dataTable[2:end,2:end]))
(n,d) = size(X)

# Standardize columns
include("misc.jl")
(X,mu,sigma) = standardizeCols(X)

# Plot matrix as image
using PyPlot
figure(1)
clf()
imshow(X)

include("PCA.jl")
model = PCA(X,2)
compressed = model.compress(X);
# Show scatterplot of 2 random features
# j1 = rand(1:d)
# j2 = rand(1:d)
figure(2)
clf()
plot(compressed[:,1],compressed[:,2],".")
for i in rand(1:n,10)
    annotate(dataTable[i+1,1],
    xy=[compressed[i,1],compressed[i,2]],
    xycoords="data")
end

```



2. Which trait of the animals has the largest influence (absolute value) on the first principal component? (Make sure not to forget the “+1” when looking for the name of the trait in the *dataTable*).

Answer: "hairless"

3. Which trait of the animals has the largest influence (absolute value) on the second principal component?

Answer: "grazer"

### 3.3 Data Compression

It is important to know how much of the information in our dataset is captured by the low-dimensional PCA representation. In class we discussed the “analysis” view that PCA maximizes the variance that is explained by the PCs, and the connection between the Frobenius norm and the variance of a centered data matrix  $X$ . Use this connection to answer the following:

1. How much of the variance is explained by our two-dimensional representation from the previous question?

Answer:  $1 - \frac{\|ZW - X\|^2}{\|x\|^2} = 1 - 0.8355 = 0.1645$

2. How many PCs are required to explain 50% of the variance in the data?

Answer: 14

Note: you can compute the Frobenius norm of a matrix using the function *norm*. Also, note that the “variance explained” formula from class assumes that  $X$  is already centered.

## 4 Very-Short Answer Questions

1. What is the difference between multi-label and multi-class classification?
2. We discussed “global” vs. “local” features for e-mail classification. What is an advantage of using global features, and what is advantage of using local features?

Answer: "Global Features are useful when we have an important email for all users. And Local features are useful for identifying emails for specific users."

3. Assuming we want to use the original features (no change of basis) in a linear model, what is an advantage of the “other” normal equations over the original normal equations?

Answer: It is faster when  $n < d$ .

4. What is the key advantage of stochastic gradient methods over gradient descent methods?

Answer: Solve the issue when "n" is really large and it will take a long time for gradient decent to take one step.

5. Which of the following step-size sequences lead to convergence of stochastic gradient to a stationary point?

(a)  $\alpha^t = 1/t^2$ .

(b)  $\alpha^t = 1/t$ .

(c)  $\alpha^t = \gamma/t$  (for  $\gamma > 0$ ).

(d)  $\alpha^t = 1/(t + \gamma)$  (for  $\gamma > 0$ ).

(e)  $\alpha^t = 1/\sqrt{t}$ .

(f)  $\alpha^t = 1$ .

Answer: b,c,d



6. In the language of loss functions and regularization, what is the difference between MLE and MAP?

Answer: MAP will give a score of the likelihood of  $w$  as prior, which is like the regularization, what gives a very small value for the  $w$  that is likely to be overfitting.

7. What is the difference between a generative model and a discriminative model?

8. In the MLE framework, what is the connection between the logistic loss and the sigmoid function in linear models.

Answer: If we use sigmoid function, we will get logistic loss after applying the MLE.

9. What is the significance of choosing  $k = 2$  for visualizing with PCA.

Answer: We can plot it using the 2 "features" to get a better visualization, than 2 original features.

10. With PCA, is it possible for the loss to increase if  $k$  is increased? Briefly justify your answer.

Answer: No. Because you can always set  $Z$  value to 0 to the  $k$  that is more than the previous one and get the same loss.

11. Why doesn't it make sense to do PCA with  $k > d$ ?

Answer: Because you will always get a perfect result.

12. In terms of the matrices associated with PCA ( $X, W, Z, \hat{X}$ ), where would an "eigenface" be stored?

Answer:  $W$

## Project Update (OPTIONAL for 340 STUDENTS)

Alongside Assignment 5, we are also asking you to submit a *project update*. It's ok if you haven't yet made much progress on your project, but we're making you submit this document as an excuse to get together with your project and plan. In particular, the project update should be a 1-page report covering the following:

1. **What has been done already:** for example, you might have already found a group of appropriate size, picked a topic, got access to appropriate data, read some related papers, or maybe even ran some preliminary experiments.
2. **What still needs to be:** what are the "action items" that need to be finished in order to complete the project. The purpose here is that you make a plan to finishing the different things that will need to come together in order to finish your project.

It's ok if the report is fairly short, and just written as series of bullet items addressing the two issues above.