15-388/688 - Practical Data Science: Nonlinear modeling, crossvalidation, and regularization

> Pat Virtue Carnegie Mellon University Spring 2022

Slide credits: CMU AI, Zico Kolter

# Outline

Example: return to peak demand prediction

Overfitting, generalization, and cross validation

Regularization

General nonlinear features

Kernels

Nonlinear classification

# Outline

#### Example: return to peak demand prediction

Overfitting, generalization, and cross validation

Regularization

General nonlinear features

Kernels

Nonlinear classification

# Peak demand vs. temperature (summer months)



#### Peak demand vs. temperature (all months)



#### Peak demand vs. temperature (all months)



# Linear regression fit



# "Non-linear" regression

Thus far, we have illustrated linear regression as "drawing a line through through the data", but this was really a function of our input features

Though it may seem limited, linear regression algorithms are quite powerful when applied to *non-linear features* of the input data, e.g.

$$x^{(i)} = \begin{bmatrix} (\text{High}_{\text{Temperature}^{(i)}})^2 \\ \text{High}_{\text{Temperature}^{(i)}} \\ 1 \end{bmatrix}$$

Same hypothesis class as before  $h_{\theta}(x) = \theta^T x$ , but now prediction will be a non-linear function of base input (e.g. a quadratic function)

Same least-squares solution  $\theta = (X^T X)^{-1} X^T y$ 



# Code for fitting polynomial

The only element we need to add to write this non-linear regression is the creation of the non-linear features

x = df\_daily.loc[:, "Temperature"] min\_x, rng\_x = (np.min(x), np.max(x) - np.min(x)) x = 2\*(x - min\_x)/rng\_x - 1.0 y = df\_daily.loc[:, "Load"] X = np.vstack([x\*\*i for i in range(poly\_degree, -1, -1)]).T theta = np.linalg.solve(X.T.dot(X), X.T.dot(y))

**Output learned function:** 

```
x0 = 2*(np.linspace(xlim[0], xlim[1],1000) - min_x)/rng_x - 1.0
X0 = np.vstack([x0**i for i in range(poly_degree,-1,-1)]).T
y0 = X0.dot(theta)
```









# Linear regression with many features

Suppose we have m examples in our data set and n = m features (plus assumption that features are linearly independent, though we'll always assume this)

Then  $X \in \mathbb{R}^{m \times n}$  is a square matrix, and least squares solution is:  $\theta = (X^T X)^{-1} X^T Y = X^{-1} X^T Y = X^{-1} Y$ 

and we therefore have  $X\theta = y$  (i.e., we fit data exactly)

Note that we can *only* perform the above operations when *X* is square, though if we have *more* features than examples, we can still get an exact fit by simply discarding features

# Outline

Example: return to peak demand prediction

#### Overfitting, generalization, and cross validation

Regularization

General nonlinear features

Kernels

Nonlinear classification

## **Generalization error**

The problem with the canonical machine learning problem is that we don't *really* care about minimizing this objective on the given data set

minimize 
$$\sum_{i=1}^{m} \ell(h_{\theta}(x^{(i)}), y^{(i)})$$

What we really care about is how well our function will generalize to *new examples* that we *didn't* use to train the system (but which are drawn from the "same distribution" as the examples we used for training)

The higher degree polynomials exhibited *overfitting*: they actually have very *low* loss on the training data, but create functions we don't expect to generalize well

# Cartoon version of overfitting

As model becomes more complex, training loss always decreases; generalization loss decreases to a point, then starts to increase





Although it is difficult to quantify the true generalization error (i.e., the error of these algorithms over the *complete* distribution of possible examples), we can approximate it by *holdout cross-validation* 

Basic idea is to split the data set into a training set and a holdout set



Train the algorithm on the training set and evaluate on the holdout set

# Cross-validation in code

```
A simple example of holdout cross-validation:
```

```
# compute a random split of the data
np.random.seed(0)
perm = np.random.permutation(len(df_daily))
idx_train = perm[:int(len(perm)*0.7)]
idx_cv = perm[int(len(perm)*0.7):]
# scale features for each split based upon training
xt = df_daily.iloc[idx_train,0]
min_xt, rng_xt = (np.min(xt), np.max(xt) - np.min(xt))
xt = 2*(xt - min_xt)/rng_xt - 1.0
xcv = 2*(df_daily.iloc[idx_cv,0] - min_xt)/rng_xt -1
yt = df_daily.iloc[idx_train,1]
ycv = df_daily.iloc[idx_cv,1]
```

```
# compute least squares solution and error on holdout and training
X = np.vstack([xt**i for i in range(poly_degree,-1,-1)]).T
theta = np.linalg.solve(X.T.dot(X), X.T.dot(yt))
err_train = 0.5*np.linalg.norm(X.dot(theta) - yt)**2/len(idx_train)
err cv = 0.5*np.linalg.norm(Xcv.dot(theta) - ycv)**2/len(idx_cv)
```

#### Parameters and hyperparameters

We refer to the  $\theta$  variables as the *parameters* of the machine learning algorithm

But there are other quantities that also affect the classifier: degree of polynomial, amount of regularization, etc; these are collectively referred to as the *hyperparameters* of the algorithm

Basic idea of cross-validation: use training set to determine the parameters, use holdout set to determine the hyperparameters

#### Illustrating cross-validation



#### Training and cross-validation loss by degree



## Training and cross-validation loss by degree



# K-fold cross-validation

A more involved (but actually slightly more common) version of cross validation

Split data set into k disjoint subsets (folds); train on k - 1 and evaluate on remaining fold; repeat k times, holding out each fold once



Report average error over all held out folds

### K-fold cross-validation

K=4-fold example:

Fold 1	Fold 2	Fold 3	Fold 4
Train	Train	Train	Val
Train	Train	Val	Train
Train	Val	Train	Train
Val	Train	Train	Train

# Variants

**Leave-one-out cross-validation:** the limit of k-fold cross-validation, where each fold is only a single example (so we are training on all other examples, testing on that one example)

[Somewhat surprisingly, for least squares this can be computed *more* efficiently than k-fold cross validation, same complexity solving for the optimal  $\theta$  using matrix equation]

**Stratified cross-validation:** keep an approximately equal percentage of positive/negative examples (or any other feature), in each fold

**Warning:** k-fold cross validation is *not* always better (e.g., in time series prediction, you would want to have holdout set all occur after training set)

# Poll 1

Say you are choosing amongst 10 different values of a polynomial degree, and you want to do K=5-fold cross-validation.

How many times do I have to train my model?

A. 5

- B. 10
- C. 15
- D. 50
- E. 5<sup>10</sup>
- F. 10<sup>5</sup>

# Outline

Example: return to peak demand prediction

Overfitting, generalization, and cross validation

#### Regularization

General nonlinear features

Kernels

Nonlinear classification



# Regularization

We have seen that the degree of the polynomial acts as a natural measure of the "complexity" of the model, higher degree polynomials are more complex (taken to the limit, we fit any finite data set exactly)

But fitting these models also requires extremely *large* coefficients on these polynomials

For 50-degree polynomial, the first few coefficients are  $\theta = -3.88 \times 10^6$ ,  $7.60 \times 10^6$ ,  $3.94 \times 10^6$ ,  $-2.60 \times 10^7$ , ...

This suggests an alternative way to control model complexity: keep the *weights small* (**regularization**)

#### **Regularized least squares**

For least squares, there is a simple solution to the regularized loss minimization problem  $\min_{\theta} \sum_{i=1}^{m} (\theta^T x^{(i)} - y^{(i)})^2 + \lambda \|\theta\|_2^2$ 

What happen as  $\lambda$  goes to zero?

What happen as  $\lambda$  gets really big?

# **Regularized loss minimization**

This leads us back to the regularized loss minimization problem we saw before, but with a bit more context now:

minimize 
$$\sum_{i=1}^{m} \ell(h_{\theta}(x^{(i)}), y^{(i)}) + \frac{\lambda}{2} \|\theta\|_{2}^{2}$$

This formulation trades off loss on the *training* set with a penalty on high values of the parameters

By varying  $\lambda$  from zero (no regularization) to infinity (infinite regularization, meaning parameters will all be zero), we can sweep out different sets of model complexity

#### 50 degree polynomial fit



#### 50 degree polynomial fit with $\lambda = 1$



#### Training/cross-validation loss by regularization


#### Training/cross-validation loss by regularization



#### **Regularized least squares**

For least squares, there is a simple solution to the regularized loss minimization problem

minimize 
$$\sum_{i=1}^{m} \left(\theta^T x^{(i)} - y^{(i)}\right)^2 + \lambda \|\theta\|_2^2$$

Taking gradients by the same rules as before gives:

$$\nabla_{\theta} \left( \sum_{i=1}^{m} \left( \theta^T x^{(i)} - y^{(i)} \right)^2 + \lambda \|\theta\|_2^2 \right) = 2X^T (X\theta - y) + 2\lambda\theta$$

Setting gradient equal to zero leads to the solution  $2X^T X \theta + 2\lambda \theta = 2X^T y \implies \theta = (X^T X + \lambda I)^{-1} X^T y$ 

Looks just like the normal equations but with an additional  $\lambda I$  term

# Outline

Example: return to peak demand prediction

Overfitting, generalization, and cross validation

Regularization

General nonlinear features

Kernels

Nonlinear classification

# More general features

We previously described polynomial features for a *single* raw input, but if our raw input is itself multi-variate, how do we define polynomial features?



# More general features

We previously described polynomial features for a *single* raw input, but if our raw input is itself multi-variate, how do we define polynomial features?

$$\begin{array}{c} \phi(x) \\ \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \rightarrow \\ \begin{bmatrix} x_1^2 \\ x_1 \\ x_2^2 \\ x_2 \\ x_1 x_2 \\ 1 \end{bmatrix}$$



# Notation for more general features

We previously described polynomial features for a *single* raw input, but if our raw input is itself multi-variate, how do we define polynomial features?

Deviating a bit from past notion, for precision here we're going to use  $x^{(i)} \in \mathbb{R}^k$  to denote the *raw* inputs, and  $\phi^{(i)} \in \mathbb{R}^n$  to denote the input features we construct (also common to use the notation  $\phi(x^{(i)})$ )

We'll also drop (*i*) superscripts, but important to understand we're transforming *each* feature this way

E.g., for the high temperature:

$$\alpha = [\text{High}_\text{Temperature}], \quad \phi = \begin{vmatrix} x^2 \\ x \end{vmatrix}$$

# Polynomial features in general

One possibility for higher degree polynomials is to just use an independent polynomial over each dimension (here of degree d)

$$x \in \mathbb{R}^{k} \Longrightarrow \phi = \begin{bmatrix} x_{1}^{d} \\ \vdots \\ x_{1} \\ \vdots \\ x_{k}^{d} \\ \vdots \\ x_{k} \\ 1 \end{bmatrix} \in \mathbb{R}^{kd+1}$$

But this ignores cross terms between different features, i.e., terms like  $x_1 x_2^2 x_k$ 

# Polynomial features in general

A better generalization of polynomials is to include all polynomial terms between raw inputs up to degree d

$$x \in \mathbb{R}^k \Longrightarrow \phi = \left\{ \prod_{i=1}^k x_i^{b_i} : \sum_{i=1}^k b_i \le d \right\} \in \mathbb{R}^{\binom{k+d}{k}}$$

Code to generate all polynomial features with degree exactly *d*:

from itertools import combinations\_with\_replacement
[np.prod(a) for a in combinations\_with\_replacement(x, d)]

Code to generate all polynomial features with degree up to d

[np.prod(a) for i in range(d+1) for a in combinations\_with\_replacement(x,i)]

# Code for general polynomials

The following code efficiently (relatively) generates all polynomials up to degree d for an entire data matrix X

It is using the same logic as above, but applying it to entire columns of the data at a time, and thus only needs one call to <code>combinations\_with\_replacement</code>

# Radial basis functions (RBFs)

Consider the following function that is related to the distance between our input x and some point  $\mu$ .

$$e^{\frac{\|x-\mu\|_2^2}{2\sigma^2}}$$

What does it look like if we plot this function versus x for  $\mu = 10$ ?

$$e^{\frac{\|x-10\|_2^2}{2\sigma^2}}$$

#### Radial basis functions (RBFs)

$$\frac{\|x - \mu\|_2^2}{2\sigma^2}$$

def rbf(x, mu, sigma):
 return np.exp(-(x-mu).T@(x-mu)/(2\*sigma\*\*2))



#### Linear combination of RBFs

$$\phi_k(x) = e^{\frac{\|x - \mu^{(k)}\|_2^2}{2\sigma^2}}$$

$$h_{\theta}(x) = \sum_{k}^{K} \theta_{k} \phi_{k}(x)$$





Linear combination of RBFs



# Radial basis functions (RBFs)

For  $x \in \mathbb{R}^k$ , select some set of *p* centers,  $\mu^{(1)}$ , ...,  $\mu^{(p)}$  (we'll discuss shortly how to select these), and create features

$$\phi = \left\{ \exp\left(-\frac{\|x - \mu^{(i)}\|_{2}^{2}}{2\sigma^{2}}\right) : i = 1, \dots, p \right\} \bigcup \{1\} \in \mathbb{R}^{p+1}$$

**Very important:** need to normalize columns of *X* (i.e., different features), to all be the same range, or distances won't be meaningful

(Hyper)parameters of the features include the choice of the p centers, and the choice of the bandwidth  $\sigma$ 

Choose centers, i.e., to be a uniform grid over input space, can choose  $\sigma$  e.g. using cross validation (don't do this, though, more on this shortly)

#### Example radial basis function

Example:

 $x = [High\_Temperature],$ 

$$\mu^{(1)} = [20], \mu^{(2)} = [25], \mu^{(3)} = [30], \mu^{(4)} = [35], \dots, \mu^{(16)} = [95], \sigma = 10$$

Leads to features:

$$\phi = \begin{bmatrix} \exp(-(\text{High}_{\text{Temperature}} - 20)^2/200) \\ \vdots \\ \exp(-(\text{High}_{\text{Temperature}} - 95)^2/200) \\ 1 \end{bmatrix}$$

#### **Example radial basis function**

Example:

 $x = [High_Temperature],$ 

 $\mu^{(1)} = [20], \mu^{(2)} = [25], \mu^{(3)} = [30], \mu^{(4)} = [35], \dots, \mu^{(16)} = [95], \sigma = 10$ 

X

Φ

	<i>x</i> <sub>1</sub>			$\phi_1$	$\phi_2$	$\phi_3$	$\phi_4$	 $\phi_{17}$
<i>x</i> <sup>(1)</sup>	25	$\phi(x)$	$\phi^{(1)}$	0.88	1.0	0.88	0.61	 1
<i>x</i> <sup>(2)</sup>	33	$\rightarrow$	$\phi^{(2)}$	0.43	0.73	0.96	0.98	 1
<i>x</i> <sup>(3)</sup>	43		$\phi^{(3)}$	0.07	0.2	0.43	0.73	 1
:			÷					

# Code for generating RBFs

The following code generates a complete set of RBF features for an entire data matrix  $X \in \mathbb{R}^{m \times k}$  and matrix of centers  $\mu \in \mathbb{R}^{p \times k}$ 

```
def rbf(X,mu,sig):
    sqdist = -2*X@mu.T + (X**2).sum(axis=1)[:,None] + (mu**2).sum(axis=1)
    return np.exp(-sqdist/(2*sig**2))
```

Important "trick" is to efficiently compute distances between *all* data points and all centers

# Difficulties with general features

The challenge with these general non-linear features is that the number of potential features grows very quickly in the dimensionality of the raw input

**Polynomials:** *k*-dimensional raw input  $\Rightarrow \binom{k+d}{k} = O(d^k)$  total features (for fixed *d*)

**RBFs:** *k*-dimensional raw input, uniform grid with *d* centers over each dimension  $\Rightarrow d^k$  total features

These quickly become impractical for large feature raw input spaces

# Practical polynomials

Don't use the full set of all polynomials, for anything but very low dimensional input data (say  $k \le 4$ )

Instead, form polynomials only of features where you know that the relationship may be important:

• E.g. Temperature<sup>2</sup> · Weekday, but not Temperature · Humidity

For binary raw inputs, no point in taking every power  $(x_i^2 = x_i)$ 

These elements do all require some insight into the problem

## **Practical RBFs**

Don't create RBF centers in a grid over your raw input space (your data will *never* cover an entire high-dimensional space, but will lie on a subset)

Instead, pick centers by randomly choosing p data points in the training set (a bit fancier, run k-means to find centers, which we'll describe later)

Don't pick  $\sigma$  using cross validation

Instead, choose the following (called the *median trick*)  $\sigma = \text{median}\left(\left\{\left\|\mu^{(i)} - \mu^{(j)}\right\|_2, i, j = 1, \dots, p\right\}\right)$ 

# Regression using training data as RBF centers

Small  $\sigma$ 

#### Medium $\sigma$

Large  $\sigma$ 



### Regression using training data as RBF centers

Small  $\sigma$ 

#### Medium $\sigma$

Large  $\sigma$ 



# Outline

Example: return to peak demand prediction

Overfitting, generalization, and cross validation

Regularization

General nonlinear features

#### Kernels

Nonlinear classification

#### Kernels

### Kernels

One of the most prominent advances in machine learning in the past 20 years (recently fallen out of favor relative to neural networks, but still can be the best-performing approach for many "medium-sized" problems)

Kernels fundamentally are about specific hypothesis function

$$h_{\theta}(x) = \sum_{i=1}^{m} \theta_i K(x, x^{(i)})$$

where  $K : \mathbb{R}^n \times \mathbb{R}^n \to \mathbb{R}$  is a kernel function

Kernels can *implicitly* represent high dimensional feature vectors *without* the need to form them explicitly (we won't prove this here, but provide a short description in the notes)

#### Kernels as high dimensional features

#### 1. Polynomial Kernel

$$K(x,z) = (1+x^T z)^d$$

is equivalent to using full degree d polynomial ( $\binom{n+d}{d}$ )-dimension) features in the raw inputs

#### 2. RBF Kernel

$$K(x,z) = \exp\left(-\frac{\|x-z\|_2^2}{2\sigma^2}\right)$$

is equivalent to a polynomial feature function with all degrees up to infinity!

## Kernels: what is the "catch"

What is the downside of using kernels?

Recall hypothesis function

$$h_{\theta}(x) = \sum_{i=1}^{m} \theta_i K(x, x^{(i)})$$

Note that we need a parameter for every training example (complexity increases with the size of the training set)

Called a *non-parametric method* (number of parameters increase with the number of data points)

Typically, complexity of resulting ML algorithm is  $O(m^2)$  (or larger), leads to impractical algorithms on large data sets

# Outline

Example: return to peak demand prediction

Overfitting, generalization, and cross validation

Regularization

General nonlinear features

Kernels

Nonlinear classification

### Nonlinear classification

Just like linear regression, the nice thing about using nonlinear features for classification is that our algorithms remain exactly the same as before

I.e., for an SVM, we just solve (using gradient descent)  

$$\min_{\theta} \sum_{i=1}^{m} \max\{1 - y^{(i)} \cdot \theta^{T} x^{(i)}, 0\} + \frac{\lambda}{2} \|\theta\|_{2}^{2}$$

Only difference is that  $x^{(i)}$  now contains non-linear functions of the input data

#### Linear SVM on cancer data set



#### Polynomial features d = 2



#### Polynomial features d = 3



#### Polynomial features d = 10



## SVM with RBF kernel $\sigma$ high



#### SVM with RBF kernel $\sigma$ medium



### SVM with RBF kernel $\sigma$ low

