

Unit 8


Support Vector Machines

EE-UY 4563/EL-GY 9143: INTRODUCTION TO MACHINE LEARNING
PROF. SUNDEEP RANGAN, WITH MODIFICATION BY YAO WANG

Learning Objectives

- ☐ Interpret weights in linear classification of images
- ☐ Describe why linear classification for images does not work
- ☐ Define the margin in linear classification
- ☐ Describe the SVM classification problem.
- ☐ Write equations for solutions of constrained optimization using the Lagrangian.
- ☐ Describe a kernel SVM problem for non-linear classification
- ☐ Implement SVM classifiers in python
- ☐ Select SVM parameters from cross-validation

Outline

- 
- Motivating example: Recognizing handwritten digits
 - Why logistic regression doesn't work well.
 - ❑ Maximum margin classifiers
 - ❑ Support vector machines
 - ❑ Kernel trick
 - ❑ Constrained optimization

MNIST Digit Classification

HANDWRITING SAMPLE FORM

NAME [REDACTED] DATE 8-3-89 CITY MINDEN CITY STATE MI ZIP 48456

This sample of handwriting is being collected for use in testing computer recognition of hand printed numbers and letters. Please print the following characters in the boxes that appear below.

0 1 2 3 4 5 6 7 8 9										0 1 2 3 4 5 6 7 8 9										0 1 2 3 4 5 6 7 8 9									
0123456789										0123456789										0123456789									
87		701		3752		80759		960941																					
158		4586		32123		832656		82																					
7481		80539		419219		67		904																					
61738		729658		75		390		5716																					

- ☐ Problem: Recognize hand-written digits
- ☐ Original problem:
 - Census forms
 - Automated processing
- ☐ Classic machine learning problem
- ☐ Benchmark

From Patrick J. Grother, NIST Special Database, 1995

A Widely-Used Benchmark

- ❑ We will look at SVM today
- ❑ Not the best algorithm
- ❑ But quite good
- ❑ ...and illustrates the main points

Classifiers [\[edit \]](#)

This is a table of some of the [machine learning](#) methods used on the database and their error rates, by type of classifier:

Type ↕	Classifier ↕	Distortion ↕	Preprocessing ↕	Error rate (%) ↕
Linear classifier	Pairwise linear classifier	None	Deskewing	7.6 ^[9]
K-Nearest Neighbors	K-NN with non-linear deformation (P2DHMDM)	None	Shiftable edges	0.52 ^[14]
Boosted Stumps	Product of stumps on Haar features	None	Haar features	0.87 ^[15]
Non-Linear Classifier	40 PCA + quadratic classifier	None	None	3.3 ^[9]
Support vector machine	Virtual SVM, deg-9 poly, 2-pixel jittered	None	Deskewing	0.56 ^[16]
Neural network	2-layer 784-800-10	None	None	1.6 ^[17]
Neural network	2-layer 784-800-10	elastic distortions	None	0.7 ^[17]
Deep neural network	6-layer 784-2500-2000-1500-1000-500-10	elastic distortions	None	0.35 ^[18]
Convolutional neural network	Committee of 35 conv. net, 1-20-P-40-P-150-10	elastic distortions	Width normalizations	0.23 ^[8]



Downloading MNIST

```
import tensorflow as tf

(Xtr,ytr),(Xts,yts) = tf.keras.datasets.mnist.load_data()

print('Xtr shape: %s' % str(Xtr.shape))
print('Xts shape: %s' % str(Xts.shape))

ntr = Xtr.shape[0]
nts = Xts.shape[0]
nrow = Xtr.shape[1]
ncol = Xtr.shape[2]
```

```
Xtr shape: (60000, 28, 28)
Xts shape: (10000, 28, 28)
```

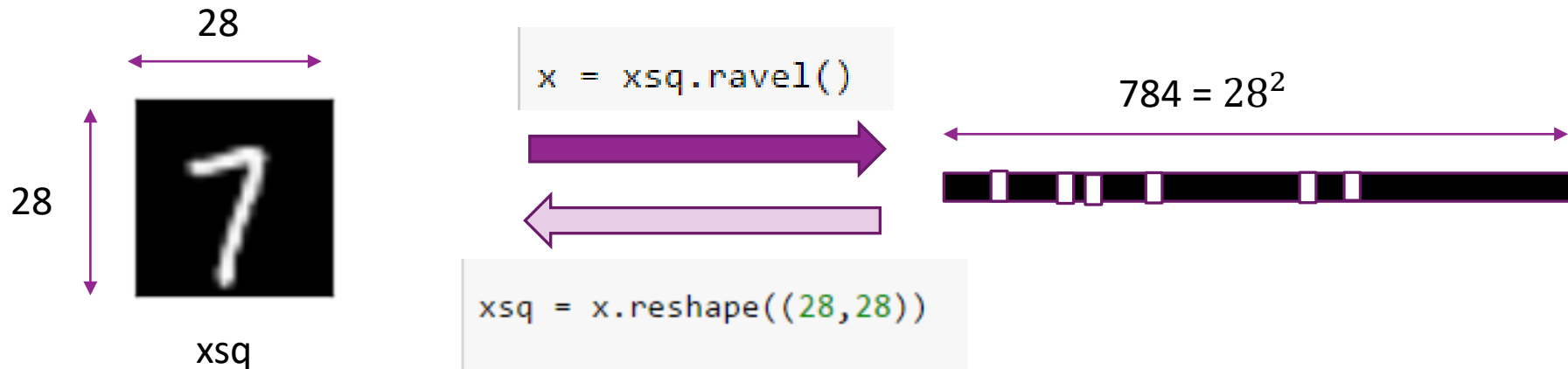
- ❑ MNIST data is available in many sources
 - Note: It has been removed from sklearn
- ❑ Tensorflow version:
 - 60000 training samples
 - 10000 test samples
- ❑ Each sample is a 28 x 28 images
- ❑ Grayscale: Pixel values $\in \{0, 1, \dots, 255\}$
 - 0 = Black and
 - 255 = White

Matrix and Vector Representation

- ❑ For this demo, we reshape data from $N \times 28 \times 28$ to $N \times 784$
- ❑ But, you can easily go back and forth
- ❑ Also, scale the pixel values from -1 to 1

```
npix = nrow*ncol
Xtr = 2*(Xtr/255 - 0.5)
Xtr = Xtr.reshape((ntr,npix))

Xts = 2*(Xts/255 - 0.5)
Xts = Xts.reshape((nts,npix))
```



$$S = \text{Mat}(x) = \begin{bmatrix} s_{11} & \cdots & s_{1,28} \\ \vdots & \vdots & \vdots \\ s_{28,1} & \cdots & s_{28,28} \end{bmatrix}$$

$$x = \text{vec}(S) = [x_1 \quad \cdots \quad x_{784}]$$

Displaying Images in Python



4 random images in the dataset

A human can classify these easily

```
def plt_digit(x):  
    nrow = 28  
    ncol = 28  
    xsq = x.reshape((nrow,ncol))  
    plt.imshow(xsq, cmap='Greys_r')  
    plt.xticks([])  
    plt.yticks([])  
  
    # Convert data to a matrix  
    X = mnist.data  
    y = mnist.target  
  
    # Select random digits  
    nplt = 4  
    nsamp = X.shape[0]  
    Iperm = np.random.permutation(nsamp)  
  
    # Plot the images using the subplot command  
    for i in range(nplt):  
        ind = Iperm[i]  
        plt.subplot(1,nplt,i+1)  
        plt_digit(X[ind,:])
```

Key command

Sample permutation is necessary for this dataset, as the original data is ordered by digits

Try a Logistic Classifier

```
ntr1 = 5000
Xtr1 = Xtr[Iperm[:ntr1],:]
ytr1 = ytr[Iperm[:ntr1]]
```

- ❑ Train on 5000 samples
 - To reduce training time.
 - In practice want to train with ~40k
- ❑ Select correct solver (lbfgs)
 - Others can be very slow. Even this will take minutes

```
from sklearn import linear_model
logreg = linear_model.LogisticRegression(verbose=10, solver='lbfgs',\
                                         multi_class='multinomial',max_iter=500)
logreg.fit(Xtr1,ytr1)
```

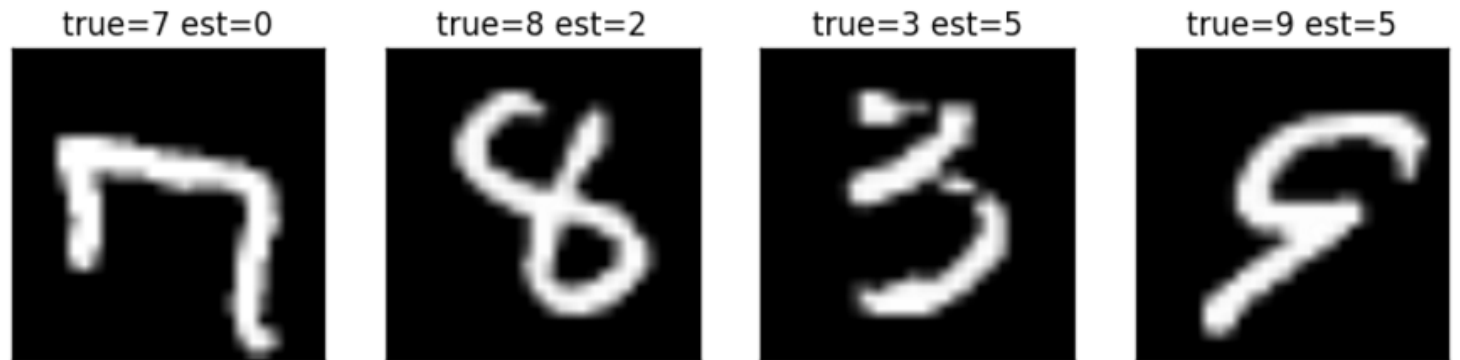
```
[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
C:\ProgramData\Anaconda3\lib\site-packages\sklearn\linear_model\logistic.py:758:
e. Increase the number of iterations.
  "of iterations.", ConvergenceWarning)
[Parallel(n_jobs=1)]: Done   1 out of   1 | elapsed:   1.3min remaining:   0.0s
[Parallel(n_jobs=1)]: Done   1 out of   1 | elapsed:   1.3min finished
```

Performance

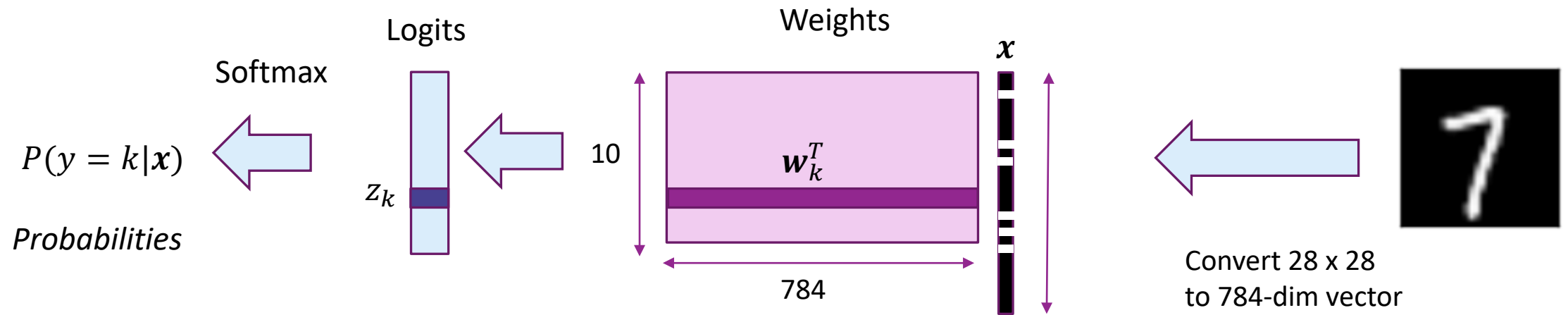
- ❑ Accuracy = 89%. Very bad
- ❑ Some of the errors seem like they should have been easy to spot
- ❑ What went wrong?

```
nts1 = 5000
Iperm_ts = np.random.permutation(nts)
Xts1 = Xts[Iperm_ts[:nts1],:]
yts1 = yts[Iperm_ts[:nts1]]
yhat = logreg.predict(Xts1)
acc = np.mean(yhat == yts1)
print('Accuracy = {0:f}'.format(acc))
```

Accuracy = 0.891000



Recap: Logistic Classifier

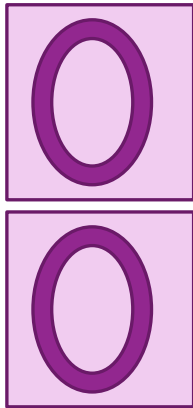


- ❑ Each **logit** $z_k = \mathbf{w}_k^T \mathbf{x}$ = inner product with **weight** \mathbf{w}_k with digit \mathbf{x} , $k = 0, \dots, 9$
- ❑ Will select $\hat{y} = \arg \max_k P(y = k|\mathbf{x}) = \arg \max_k z_k$
 - Output z_k which is largest
- ❑ When is z_k large?

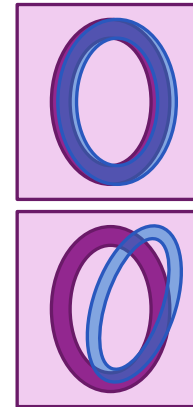
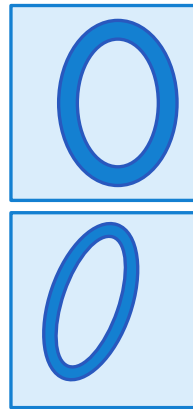
Interpreting the Logistic Classifier Weights

- A logit $z_k = \mathbf{w}_k^T \mathbf{x}$ is high when there is high **overlap** between \mathbf{w}_k with digit x
 - Visualize each weight as an image
 - Suppose pixels are 0 or 1
 - $z_k = \mathbf{w}_k^T \mathbf{x} = \sum_i w_{ki} x_i =$ number of pixels that overlap with \mathbf{w}_k and \mathbf{x}
- **Conclusion:** Small variations in digits can cause low overlap

Weight for digit "0"
 \mathbf{w}_0



Digit
 x



Overlap **high** $\Rightarrow x$ *is* digit 0

Overlap **low** $\Rightarrow x$ *not* digit 0

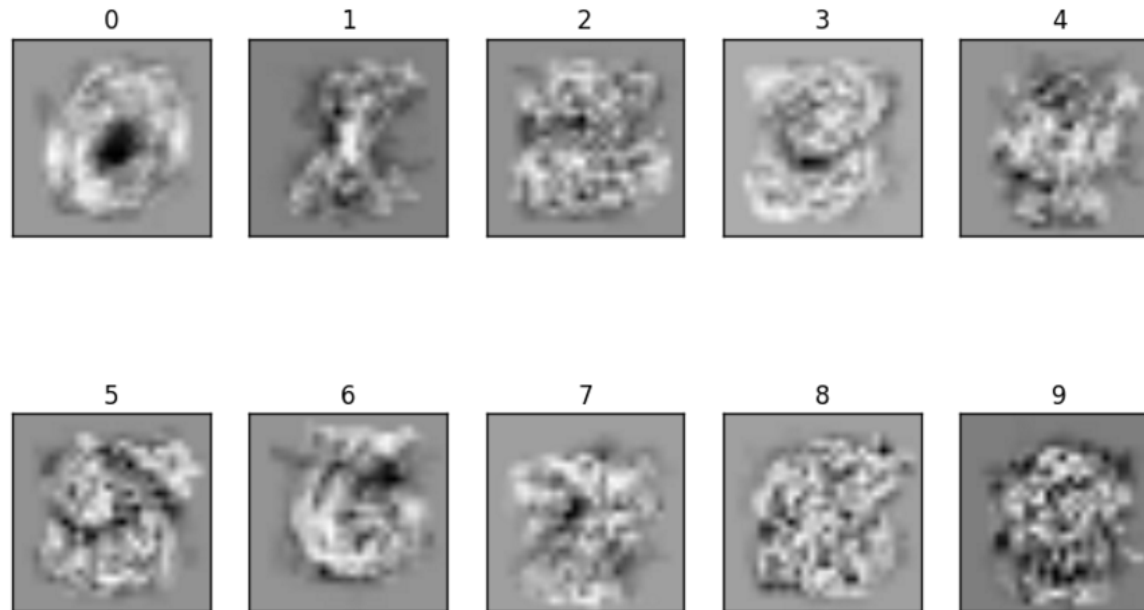
Example with Actual Digits

- Take weight w from a random digit “2”
- Inner products $z = w^T x$ are only slightly higher for other digits “2”
- Cannot tell which digit is correct from the inner product $z = w^T x$



Visualizing the Weights


- ❑ Optimized weights of the classifier
- ❑ Blurry versions of image to try to capture rotations, translations, ...



Problems with Logistic Classifier

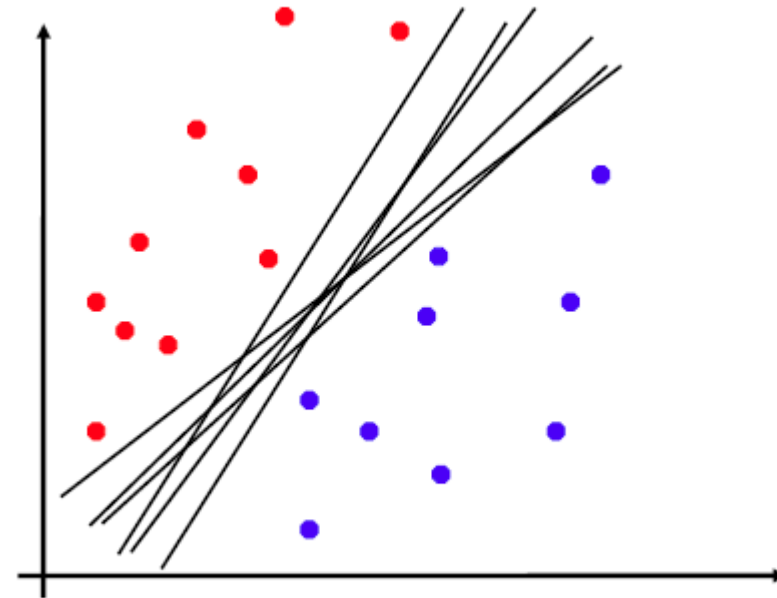
- ❑ Linear weighting cannot capture many deformities in image
 - Rotations
 - Translations
 - Variations in relative size of digit components
- ❑ Can be improved with preprocessing
 - E.g. deskewing, contrast normalization, many methods
- ❑ Is there a better classifier?

Outline

- ❑ Motivating example: Recognizing handwritten digits
 - Why logistic regression doesn't work well.
- ❑ Maximum margin classifiers
- ❑ Support vector machines
- ❑ Kernel trick
- ❑ Constrained optimization

Non-Uniqueness of Separating Plane

- Linearly separable data:
 - Can find a separating hyper-plane as a linear classifier.
- Separating hyper-plane is not unique
 - Fig. on right: Many separating planes
- Which one is optimal?



Hyperplane Basics

□ **Linear function:** $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b, \mathbf{x} \in R^d$

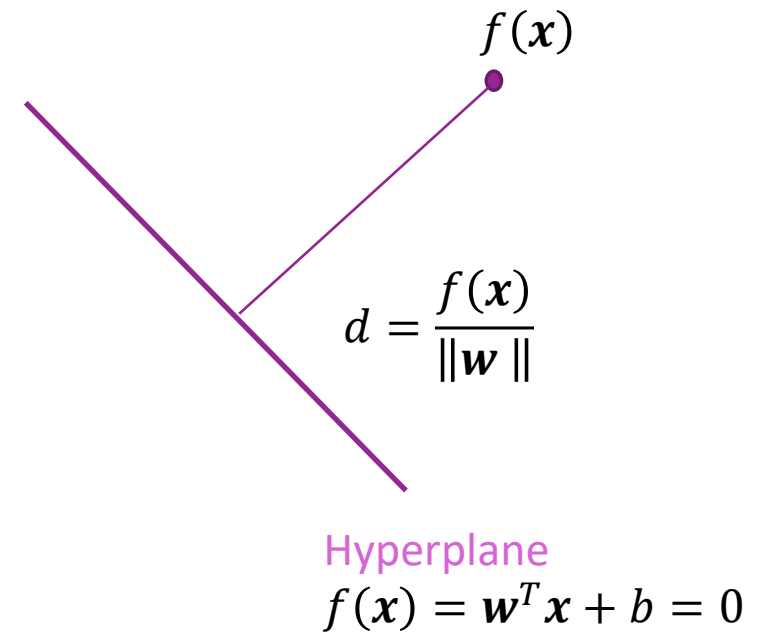
□ **Hyperplane** in d-dimensional: $f(\mathbf{x}) = 0$

□ **Parameters:**

- Weight \mathbf{w} and bias b
- Unique up to scaling:
- (b, \mathbf{w}) and $(\alpha b, \alpha \mathbf{w})$ define the same plane.
- For unique definition, we can require $\|\mathbf{w}\|=1$.

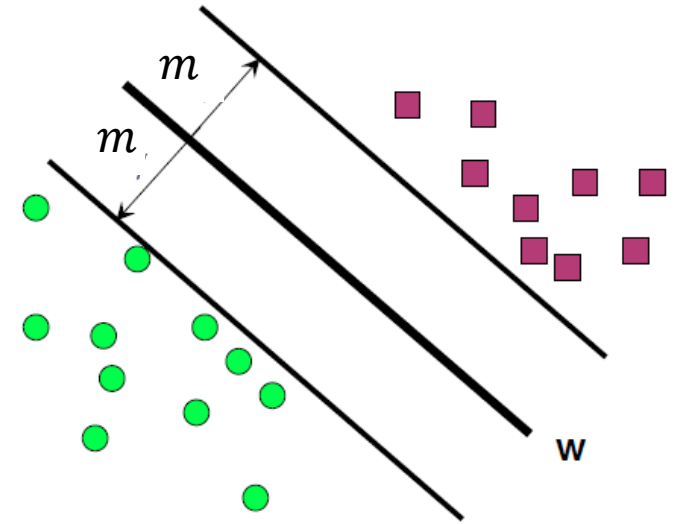
□ **Distance** of any point \mathbf{x} to the hyperplane:

- $d = f(\mathbf{x})/\|\mathbf{w}\|$, where $f(\mathbf{x}) = b + \mathbf{w}^T \mathbf{x}$.
- See ESL Sec. 4.5.
- ESL: Hastie, Tibshirani, Friedman, “The Elements of Statistical Learning”. 2nd Ed. Springer.

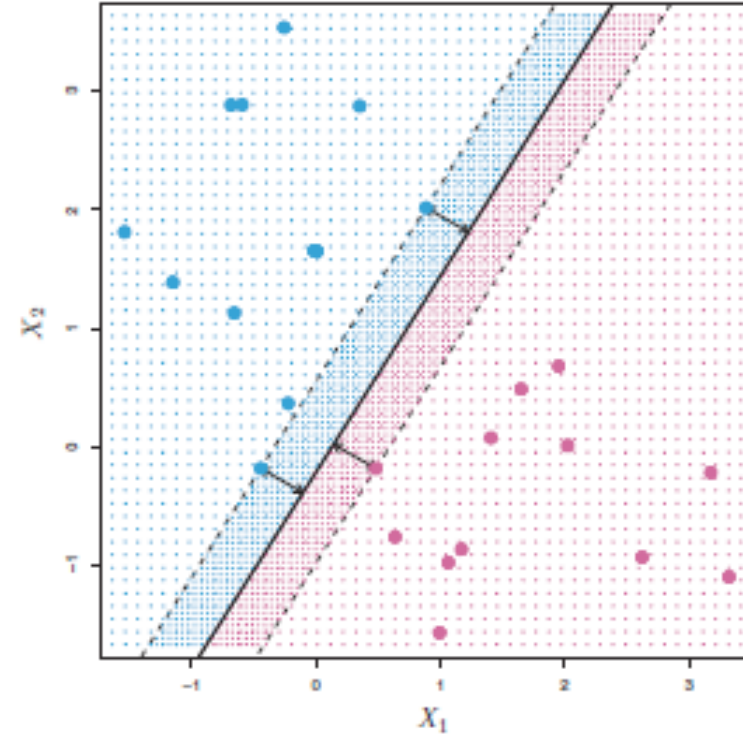
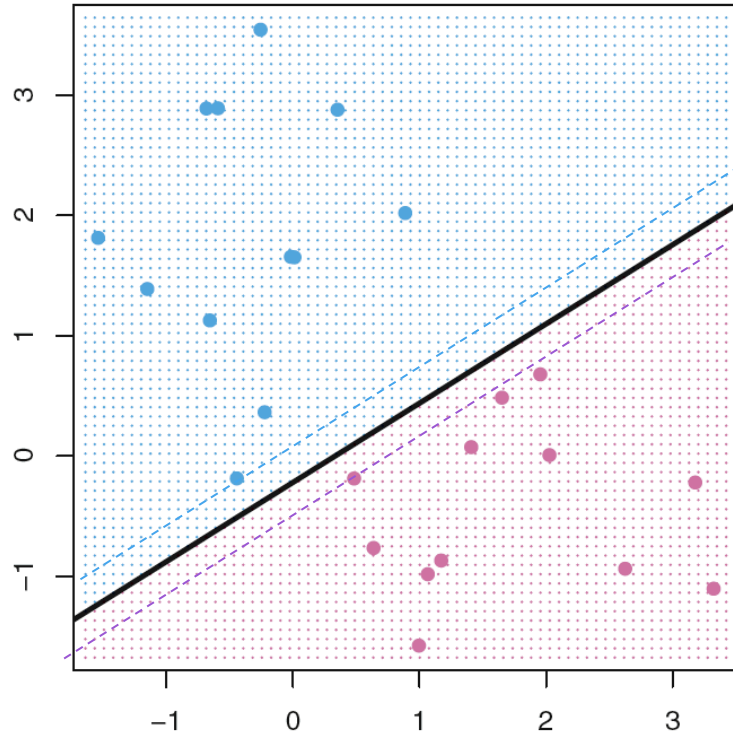


Linear Separability and Margin

- Given training data $(x_i, y_i), i = 1, \dots, N$
 - Binary class label: $y_i = \pm 1$
- Suppose it is separable with parameters (w, b)
- There must exist a $\gamma > 0$ s.t.:
 - $b + w_1x_{i1} + \dots w_dx_{id} > \gamma$ when $y_i = 1$
 - $b + w_1x_{i1} + \dots w_dx_{id} < -\gamma$ when $y_i = -1$
- Single equation form:
$$y_i(b + w_1x_{i1} + \dots w_dx_{id}) > \gamma \text{ for all } i = 1, \dots, N$$
- **Margin:** $m = \frac{\gamma}{\|w\|}$: minimal distance of a sample to the plane
 - γ is the minimum value satisfying the above constraints



Which separating plane is better ?



From Fig. 9.2 and Fig. 9.3 in ISL.

Maximum Margin Classifier

❑ For the classifier to be more robust to noise, we want to maximize the margin!

❑ Define **maximum margin classifier**

$$\max_{w, \gamma} \gamma$$

◦ Such that $y_i(b + \mathbf{w}^T \mathbf{x}) \geq \gamma$ for all i

◦ $\sum_{j=1}^d w_j^2 \leq 1$

← Maximizes the margin

← Ensures all points are correctly classified

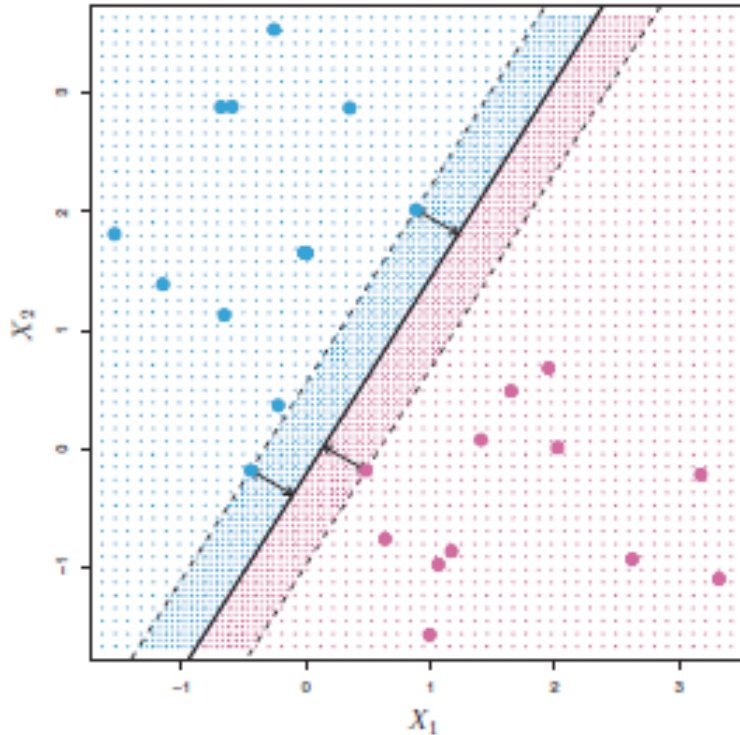
← Scaling on weights

❑ Called a **constrained optimization**

- Objective function and constraints
- More on this later.

❑ See closed form solution in Sec. 4.5.2 in ESL. Note notation difference.

Visualizing Maximum Margin Classifier



- Fig. 9.3 of ISL
- Margin determined by closest points to the line
 - The maximal margin hyperplane represents the mid-line of the **widest “slab”** that we can insert between two classes
- In this figure, there are 3 points at the margin

ISL: James, Witten, Hastie, Tibshirani, An Introduction to Statistical Learning, Springer. 2013.

Problems with MM classifier

- ❑ Data is often not perfectly separable
 - Only want to correctly separate most points

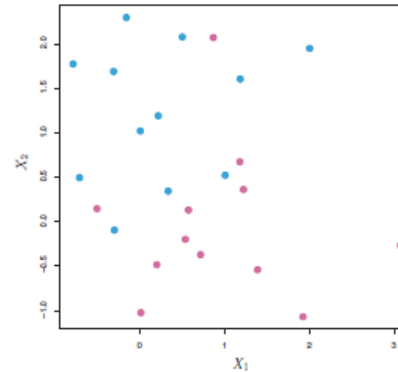


Fig. 9.4

- ❑ MM classifier is not robust
 - A single sample can radically change line

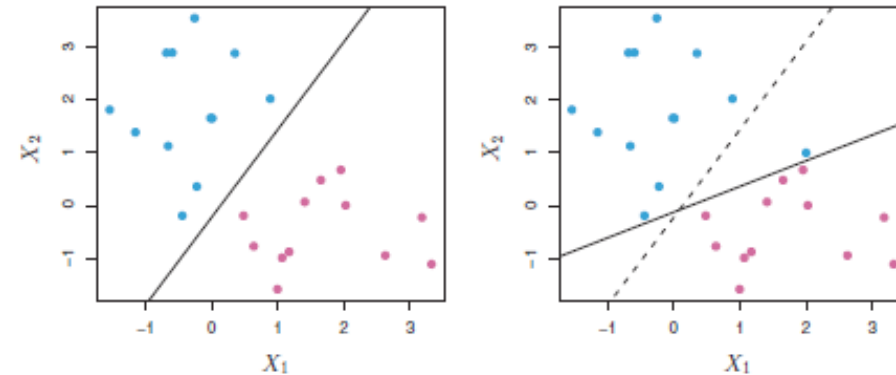


Fig. 9.5

In-Class Exercise

Found in github site: `svm_inclass.ipynb`


Problem 1. Margin

For the points below with binary labels:

- Create a scatter plot of the points with different markers for the two classes
- Find the weight and bias of the classifier that separates the two classes
- Compute the distance to the classifier boundary for the points
- Find the margin

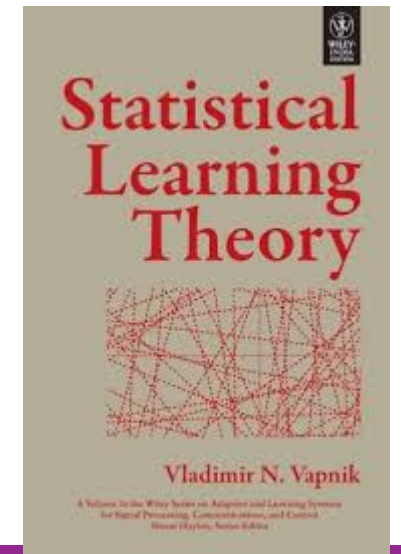
```
: ▶ 1 X = np.array([[0.5,0.5], [1,0.5],[0.5,1.75], [0.75,2.75], [1.1,2.2], [2,1], [3,1.5]])  
    2 y = np.array([1,1,1,0,0,0,0])  
    3
```


Outline

- ❑ Motivating example: Recognizing handwritten digits
 - Why logistic regression doesn't work well.
- ❑ Maximum margin classifiers
- ❑ Support vector machines
- ❑ Kernel trick
- ❑ Constrained optimization

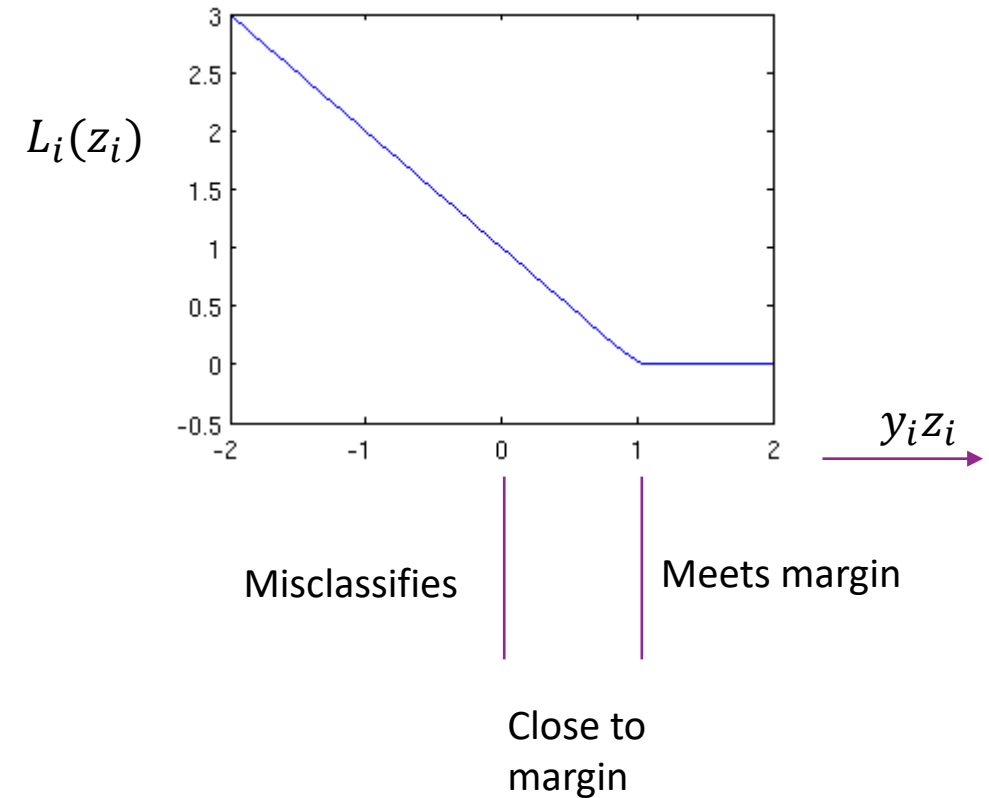
Support Vector Machine

- ❑ Support Vector Machine (SVM)
 - Vladimir Vapnik, 1963
 - But became widely-used with kernel trick, 1993
 - More on this later
- ❑ Got best results on character recognition
- ❑ Key idea: Allow “slack” in the classification
 - Support vector classifier (SVC): Directly use raw features. Good when the original feature space is roughly linearly separable
 - Support vector machine (SVM): Map the raw features to some other domain through a kernel function



Hinge Loss

- Fix $\gamma = 1$
- Want ideally: $y_i(\mathbf{w}^T \mathbf{x} + b) \geq 1$ for all samples i
 - Equivalently, $y_i z_i \geq 1$, $z_i = b + \mathbf{w}^T \mathbf{x}$
- But perfect separation may not be possible
- Define **hinge loss** or **soft margin**:
 - $L_i(\mathbf{w}, b) = \max(0, 1 - y_i z_i)$
- Starts to increase as sample is misclassified:
 - $y_i z_i \geq 1 \Rightarrow$ Sample meets margin target, $L_i(\mathbf{w}) = 0$
 - $y_i z_i \in [0, 1) \Rightarrow$ Sample margin too small, small loss
 - $y_i z_i \leq 0 \Rightarrow$ Sample misclassified, large loss



SVM Optimization

□ Given data (\mathbf{x}_i, y_i)

□ Optimization $\min_{w,b} J(\mathbf{w}, b)$

$$J(\mathbf{w}, b) = C \sum_{i=1}^N \max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b)) + \frac{1}{2} \|\mathbf{w}\|^2$$

C controls final margin

Hinge loss term
Attempts to reduce
Misclassifications

margin = $1/\|\mathbf{w}\|$

□ Constant $C > 0$ will be discussed below

□ Note: ISL book uses different naming conventions.

- We have followed convention in sklearn

Alternate Form of SVM Optimization

□ Equivalent optimization:

$$\min J_1(\mathbf{w}, b, \boldsymbol{\epsilon}), \quad J_1(\mathbf{w}, b, \boldsymbol{\epsilon}) = C \sum_{i=1}^N \epsilon_i + \frac{1}{2} \|\mathbf{w}\|^2$$

□ Subject to constraints:

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \epsilon_i \text{ for all } i = 1, \dots, N$$

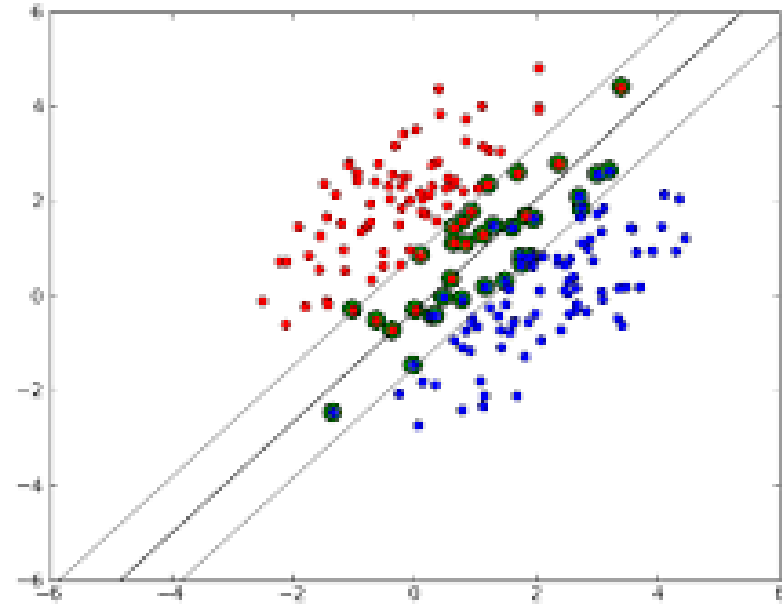
- ϵ_i = amount sample i misses margin target

□ Sometimes write as $J_1(\mathbf{w}, b, \boldsymbol{\epsilon}) = C \|\boldsymbol{\epsilon}\|_1 + \frac{1}{2} \|\mathbf{w}\|^2$

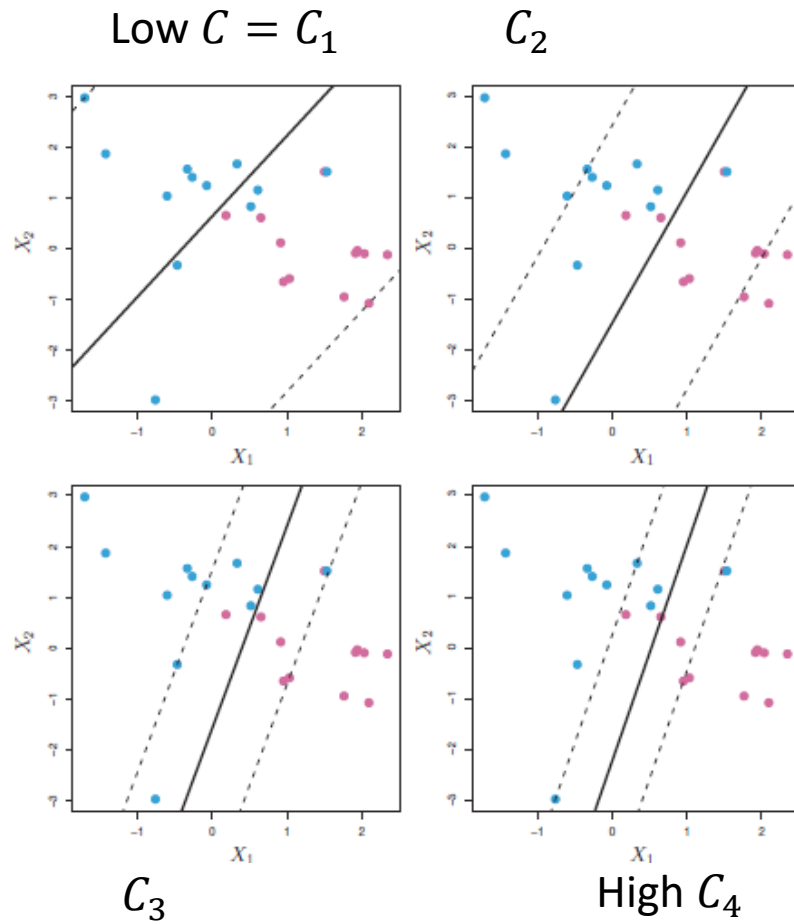
- $\|\boldsymbol{\epsilon}\|_1 = \sum_{i=1}^N \epsilon_i$ called the “one-norm”
- Generally one-norm would have absolute sign over ϵ_i .
- But in this case, when the constraint is met, $\epsilon_i \geq 0$.

Support Vectors

- ❑ **Support vectors:** Samples that either:
 - Are exactly on margin: $y_i(\mathbf{w}^T \mathbf{x}_i + b) = 1$
 - Or, on wrong side of margin: $y_i(\mathbf{w}^T \mathbf{x}_i + b) \leq 1$
- ❑ Changing samples that are not SVs
 - Does not change solution
 - Provides robustness



Illustrating Effect of C



□ Fig. 9.7 of ISL

- Note: C has opposite meaning in ISL than python
- Here, we use python meaning

□ Low C :

- Leads to large margin
- But allow many violations of margin.
- Many more SVs
- Reduces variance by using more samples

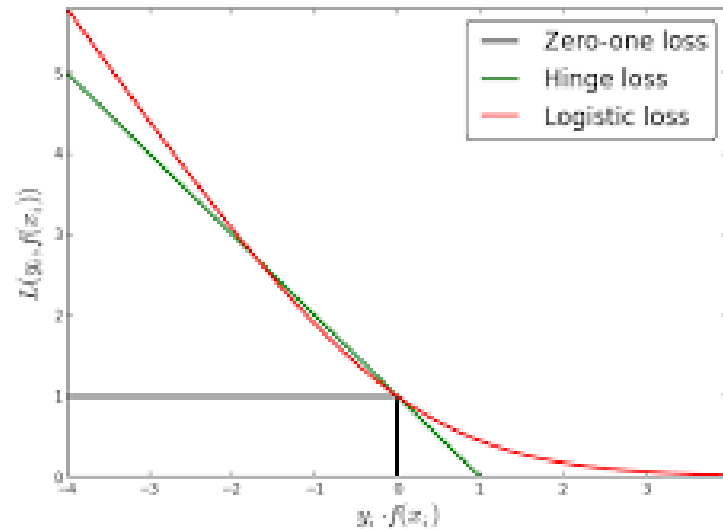
□ Large C :

- Leads to small margin
- Reduce number of violations, and fewer SVs.
- Highly fit to data. Low bias, higher variance
- More chance to overfit

Relation to Logistic Regression

□ Logistic regression also minimizes a loss function:

$$J(\mathbf{w}, b) = \sum_{i=1}^N L_i(\mathbf{w}, b), \quad L_i(\mathbf{w}, b) = \ln P(y_i | \mathbf{x}_i) = -\ln(1 + e^{-y_i z_i})$$



In-Class Exercise

Problem 2. Minimizing the Hinge Loss

For the data below, first create a scatter plot of the points with different markers for the two classes. You should see that the data is not linearly separable.

Then, consider a set of classifiers:

$$\hat{y} = \text{sign}(z), \quad z = w \cdot x + b$$


Use the w below, plot the hinge loss as a function of the bias b where the hinge loss is:

$$J = \sum \max(0, 1 - y_i w \cdot x_i)$$

Here $y_i = 2y_i - 1$ so that it is a value +1 or -1. Find the b that minimizes the hinge loss and plot the boundary of the classifier.

```
1 X = np.array([[0.5,0.5], [1,0.5],[0.5,1.75], [2,2], [0.75,0.75], [0.75,2.75], [1.1,2.2], [2,1], [3,1.5]])
2 y = np.array([1,1,1,1,0,0,0,0,0])
3
4 w = np.array([1.5, 1])
5 w = w / np.linalg.norm(w)
```

Outline

- ❑ Motivating example: Recognizing handwritten digits
 - Why logistic regression doesn't work well.
- ❑ Maximum margin classifiers
- ❑ Support vector machines
- ❑ Kernel trick
- ❑ Constrained optimization

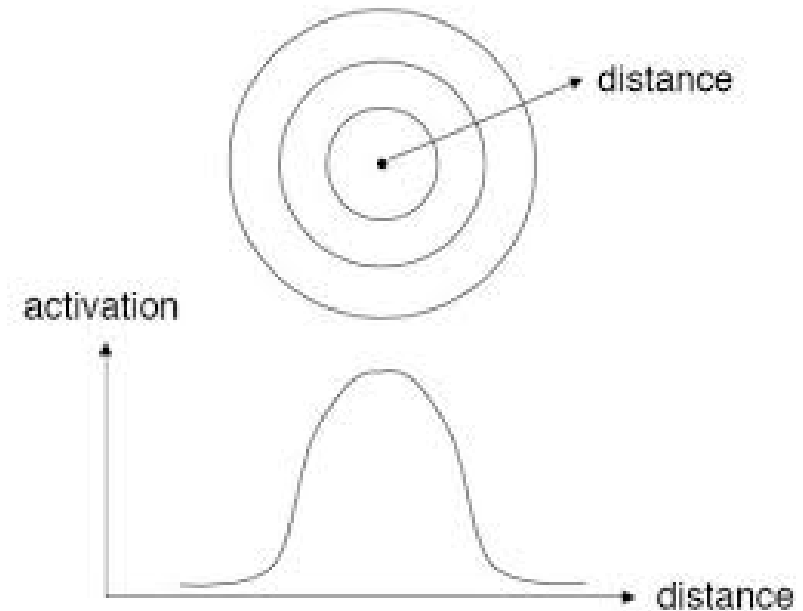
The Kernel Function

□ Kernel function:

- Function $K(x_i, x)$
- Key function for SVMs and kernel classifiers
- Measures “similarity” between new sample x and training sample x_i

□ Typical property

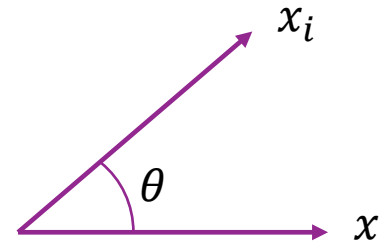
- x_i, x close $\Rightarrow K(x_i, x)$ maximum value
- x_i, x far $\Rightarrow K(x_i, x) \approx 0$



Common Kernels

Linear SVM:

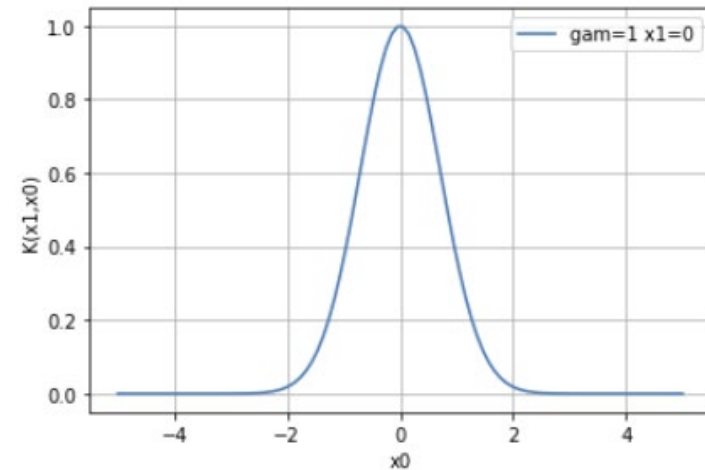
- $K(x_i, x) = x_i^T x = \|x_i\| \|x\| \cos \theta$
- Maximum when angle between vectors is small



Radial basis function:

$$K(x_i, x) = \exp[-\gamma \|x - x_i\|^2]$$

- $1/\gamma$ indicates width of kernel



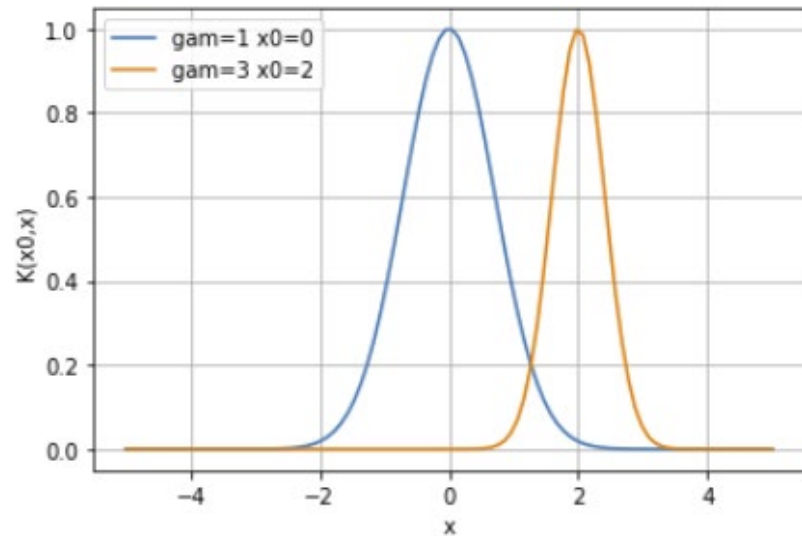
Polynomial kernel: $K(x_i, x) = |x_i^T x|^d$

- Typically $d=2$

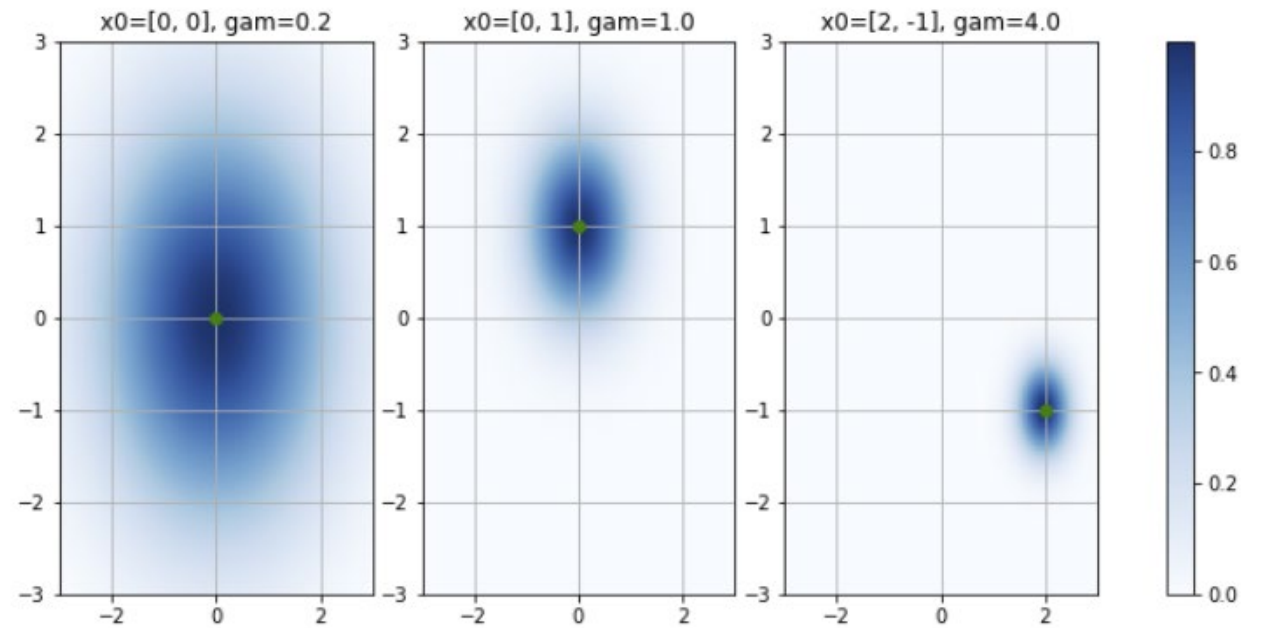
RBF Kernel Examples

□ RBF kernel: $K(x_0, x) = \exp[-\gamma \|x - x_0\|^2]$

- Peak value of 1 at $x = x_0$
- Width $\propto \frac{1}{\gamma}$



RBFs in 1D



RBFs in 2D

Kernel Classifier

□ Given:

- Training data (x_i, y_i) with binary labels $y_i = \pm 1$
- Kernel $K(x_i, x)$

□ To classify a new point x :

- Decision function: $z = \sum_{i=1}^n y_i K(x_i, x)$
- Classify: $\hat{y} = \text{sign}(z)$

□ Idea:

- z is large positive when x is close to samples x_i with $y_i = 1$
- z is large negative when x is close to samples x_i with $y_i = -1$

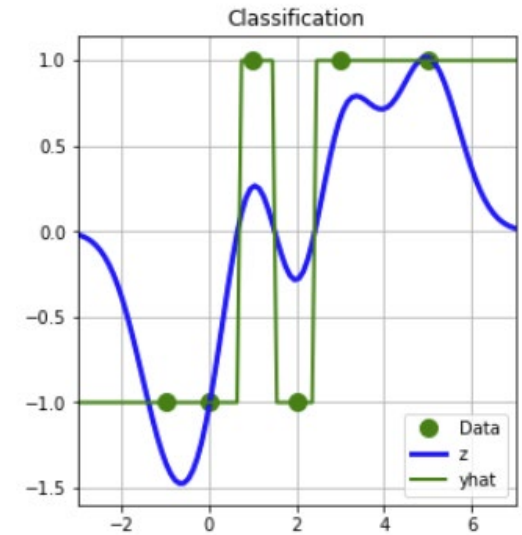
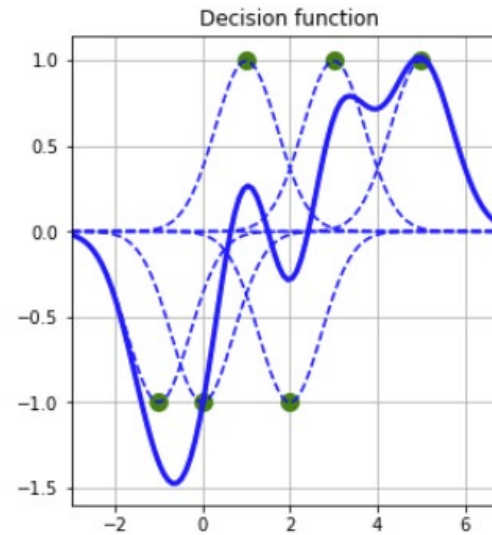
□ Kernel classifiers are a subject on their own

- We just mention them here to explain connection to SVMs

Example in 1D

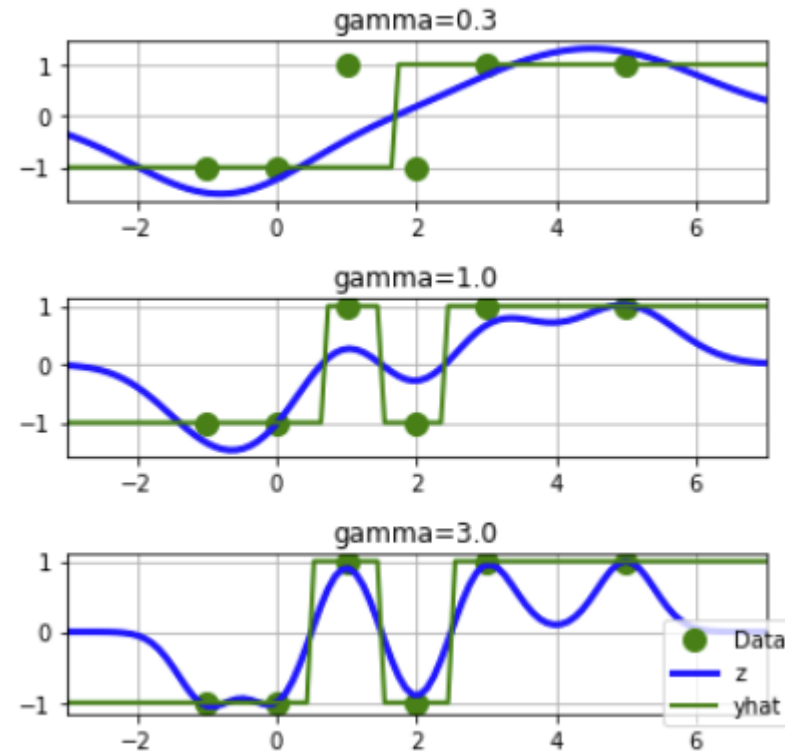
- Example data with 6 points (x_i, y_i)
 - RBF kernel: $K(x_i, x) = e^{-\gamma(x_i - x)^2}$, $\gamma = 1$
- Decision function:
 - $z = \sum_{i=1}^n y_i K(x_i, x)$
 - Sum of bell curves
 - Positive when near positive samples
 - Negative when near negative samples
- Classification:
 - $\hat{y} = \text{sign}(z)$

i	1	2	3	4	5	6
x_i	-1	0	1	2	3	5
y_i	-1	-1	1	-1	1	1



Effect of Gamma

- ❑ Same data as before
- ❑ RBF kernel: $K(x_i, x) = e^{-\gamma(x_i - x)^2}$
- ❑ As γ increases:
 - Decision function $z \approx y_i$ when $x = x_i$
 - Classifier fits training data better
 - Classification region more complex
- ❑ As a classifier, higher γ results in:
 - Lower bias error
 - But, higher variance error



SVMs with Non-Linear Transformations

❑ Non-linear transformation:

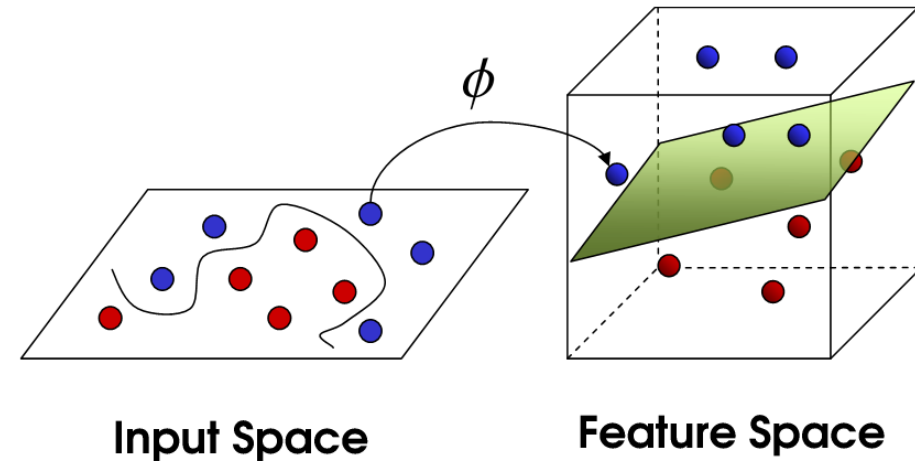
- Replace x with $\phi(x)$
- Enables more rich, non-linear classifiers
- Examples: polynomial classification

$$\phi(x) = [1, x, x^2, \dots, x^{d-1}]$$

❑ Tries to find separation in a **feature** space

❑ **Kernel trick** in SVMs:

- Makes applying non-linear transformations easy



SVM with the Transformation

- ❑ Consider SVM model with \mathbf{x} replaced by $\phi(\mathbf{x})$
- ❑ Minimize SVM cost function as before (i.e. Hinge loss + inverse margin)
- ❑ **Theorem:** The optimal weight is of the form:

$$\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \phi(\mathbf{x}_i)$$

- $\alpha_i \geq 0$ for all i
- $\alpha_i > 0$ if and only if sample i is a support vector
- Will show this fact later using results in constrained optimization

- ❑ **Consequence:** The linear discriminant on any other sample \mathbf{x} is:

$$z = b + \mathbf{w}^T \phi(\mathbf{x}) = b + \sum_{i=1}^N \alpha_i y_i \phi(\mathbf{x}_i)^T \phi(\mathbf{x})$$

————— $K(\mathbf{x}_i, \mathbf{x}) = \text{“kernel”}$

◦

Kernel Form of the SVM Classifier

□ SVM classifier can be written with the kernel $K(\mathbf{x}_i, \mathbf{x})$ and values $\alpha_i \geq 0$:

$$z = b + \sum_{i=1}^N \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}),$$

Decision function

$$\hat{y} = \text{sign}(z) = \begin{cases} 1 & \text{if } z > 0 \\ -1 & \text{if } z < 0 \end{cases}$$

Classification decision

□ **Key point:** SVM classifier is approximately Kernel classifier

□ But there are two differences:

- Weights $\alpha_i \geq 0$ on the samples (the weights are only non-zero on the SVs)
- A bias term b (can be positive or negative)

“Kernel Trick” and Dual Parameterization

□ Kernel form of SVM classifier (previous slide):

$$z = b + \sum_{i=1}^N \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}),$$
$$\hat{y} = \text{sign}(z)$$

□ Dual parameters: $\alpha_i \geq 0, i = 1, \dots, N$

- Called the dual parameters due to constrained optimization – see next section

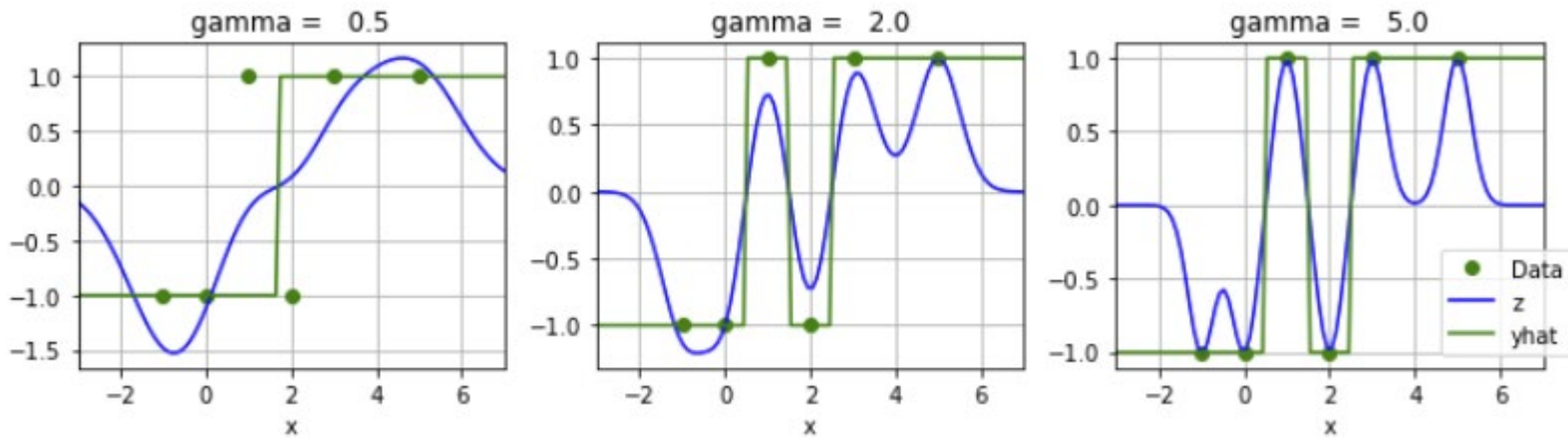
□ Kernel trick:

- Directly solve the parameters α instead of the weights \mathbf{w}
- Can show that the optimization only needs the kernel $K(\mathbf{x}_i, \mathbf{x})$
- Does not need to explicitly use $\phi(\mathbf{x})$

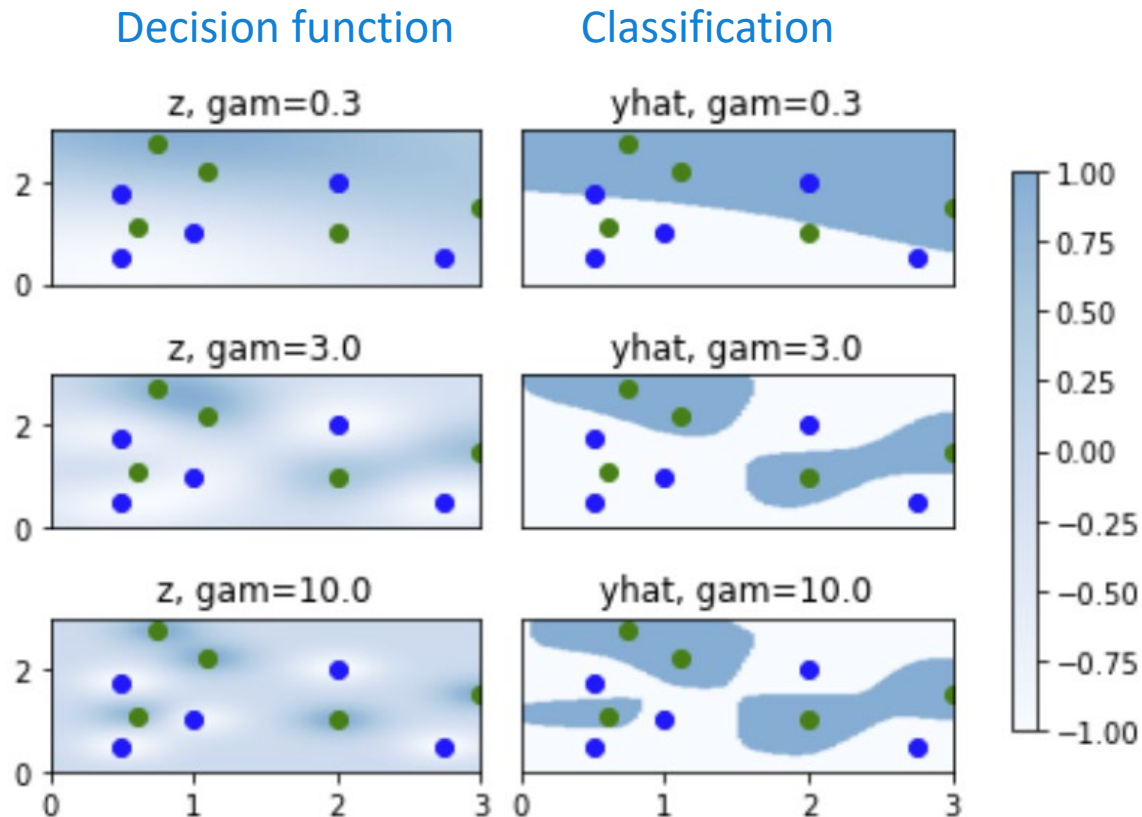
SVM Example in 1D

- ❑ Same data as in the Kernel classifier example
- ❑ Fit SVM with RBF with different γ
- ❑ Similar trends as kernel classifier: As γ increases
 - z “fits” data (x_i, y_i) closer
 - Leads to more complex decision regions.
 - Enables nonlinear decision regions

i	1	2	3	4	5	6
x_i	-1	0	1	2	3	5
y_i	-1	-1	1	-1	1	1



Example in 2D



Example:

- 10 data points with binary labels
- Fit SVM with $C = 1$ and RBF
- $\gamma = 0.3, 3$ and 10

Plot:

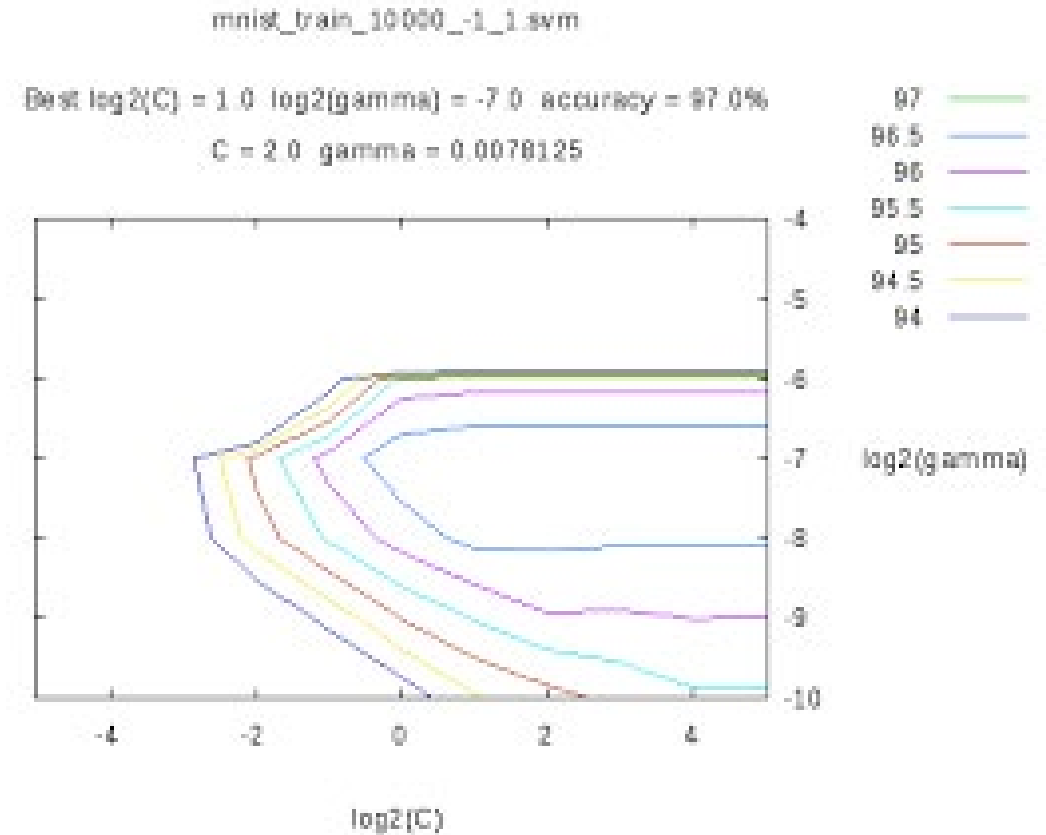
- z = linear discriminant
- $\hat{y} = \text{sign}(z)$ = classification decision

Observe: As γ increases

- Fits training data better
- More complex decision region

Parameter Selection

- ❑ For SVMs with RBFs we need to select:
 - Parameter $C > 0$ in the loss function
 - Kernel width $\gamma > 0$
- ❑ Higher C or γ
 - Fewer SVs
 - Classifiers averages over smaller set
 - Lower bias, but higher variance
- ❑ Typically select via cross-validation
 - Try out different (C, γ) pairs
 - Find which one provides highest accuracy on test set
- ❑ Python can automatically do grid search



<http://peekaboo-vision.blogspot.com/2010/09/mnist-for-ever.html>

Multi-Class SVMs

- ❑ Suppose there are K classes
- ❑ One-vs-one:
 - Train $\binom{K}{2}$ SVMs for each pair of classes
 - Test sample assigned to class that wins “majority of votes”
 - Best results but very slow
- ❑ One-vs-rest:
 - Train K SVMs: train each class k against all other classes
 - Pick class with highest z_k
- ❑ Sklearn has both options

MNIST Results

- ❑ Run classifier
- ❑ Very slow
 - Several minutes for 40,000 samples
 - Slow in training and test
 - Major drawback of SVM
- ❑ Accuracy ≈ 0.984
 - Much better than logistic regression
- ❑ Can get better with:
 - pre-processing
 - More training data
 - Optimal parameter selection

```
from sklearn import svm
```

```
# Create a classifier: a support vector classifier
```

```
svc = svm.SVC(probability=False, kernel="rbf", C=2.8, gamma=.0073, verbose=10)
```

```
svc.fit(Xtr, ytr)
```

```
[LibSVM]
```

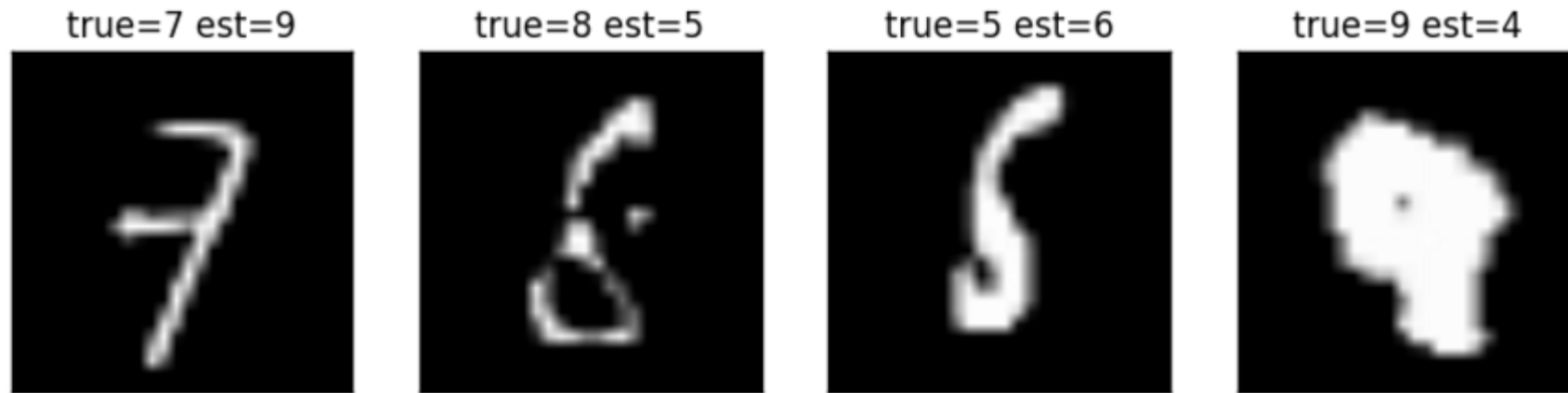
```
SVC(C=2.8, cache_size=200, class_weight=None, coef0=0.0,  
    decision_function_shape=None, degree=3, gamma=0.0073, kernel='rbf',  
    max_iter=-1, probability=False, random_state=None, shrinking=True,  
    tol=0.001, verbose=10)
```

```
yhat1 = svc.predict(Xts)  
acc = np.mean(yhat1 == yts)  
print('Accuracy = {0:f}'.format(acc))
```


Accuracy = 0.984000

MNIST Errors

- Some of the error are hard even for a human



Outline

- ❑ Motivating example: Recognizing handwritten digits
 - Why logistic regression doesn't work well.
- ❑ Maximum margin classifiers
- ❑ Support vector machines
- ❑ Kernel trick
- ❑ Constrained optimization

Constrained Optimization

- ❑ In many problems, variables are constrained
- ❑ Constrained optimization formulation:
 - **Objective:** Minimize $f(\mathbf{w})$
 - **Constraints:** $g_1(\mathbf{w}) \leq 0, \dots, g_M(\mathbf{w}) \leq 0$
- ❑ Examples:
 - Minimize the mpg of a car subject to a cost or meeting some performance
 - In ML: weight vector may have constraints from physical knowledge
- ❑ Often write constraints in vector form: Write $g(\mathbf{w}) \leq 0$

$$g(\mathbf{w}) = [g_1(\mathbf{w}), \dots, g_m(\mathbf{w})]^T$$

Lagrangian

□ Constrained optimization: $\text{Min } f(\mathbf{w}) \text{ s.t. } g(\mathbf{w}) \leq 0$

□ Consider first a single constraint: $g(\mathbf{w})$ is a scalar

□ Define **Lagrangian**: $L(\mathbf{w}, \lambda) = f(\mathbf{w}) + \lambda g(\mathbf{w})$

- \mathbf{w} is called the **primal** variable
- λ is called the **dual** variable

□ **Dual minimization**: Given a dual parameter λ , minimize

$$\hat{\mathbf{w}}(\lambda) = \arg \min_{\mathbf{w}} L(\mathbf{w}, \lambda), \quad L^*(\lambda) = \min_{\mathbf{w}} L(\mathbf{w}, \lambda)$$

- Minimizes a weighted combination of objective and constraint.
- Higher $\lambda \Rightarrow$ Weight constraint more (try to make $g(\mathbf{w})$ smaller)
- Lower $\lambda \Rightarrow$ Weight objective more (try to make $f(\mathbf{w})$ smaller)

KKT Conditions

□ Given objective $f(\mathbf{w})$ and constraint $g(\mathbf{w})$

□ **KKT Conditions:** $\hat{\mathbf{w}}, \hat{\lambda}$ satisfy:

- $\hat{\mathbf{w}}$ minimizes the Lagrangian: $\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} L(\mathbf{w}, \hat{\lambda})$
- Either
 - $g(\hat{\mathbf{w}}) = 0$ and $\hat{\lambda} \geq 0$ [active constraint]
 - $g(\hat{\mathbf{w}}) < 0$ and $\hat{\lambda} = 0$ [inactive constraint]

□ **Theorem:** Under some technical conditions,

- if $\hat{\mathbf{w}}, \hat{\lambda}$ are local minima of the constrained optimization, they must satisfy KKT conditions

General Procedure for Single Constraint

□ Suppose:

- $\mathbf{w} = (w_1, \dots, w_d)^T$: d unknown primal variables
- $g(\mathbf{w}) \leq 0$: scalar constraint

□ Case 1: Assume constraint is active:

- Solve \mathbf{w} and λ : $\partial L(\mathbf{w}, \lambda) / \partial w_j = 0$ and $g(\mathbf{w}) = 0$ (resulting from setting $\partial L(\mathbf{w}, \lambda) / \partial \lambda = 0$)
- $d + 1$ unknowns and $d + 1$ equations
- Verify that $\lambda \geq 0$

□ Case 2: Assume constraint is inactive

- Solve primal objective $\partial f(\mathbf{w}) / \partial w_j = 0$ ignoring constraint
- d unknowns and d equations
- Verify that constraint is satisfied: $g(\mathbf{w}) \leq 0$

KKT Conditions Illustrated

- Example 1: Constraint is “active”

$$\min_w w^2 \quad s.t. \quad w + 1 \leq 0$$

- Example 2: Constraint is “inactive”

$$\min_w w^2 \quad s.t. \quad w - 1 \leq 0$$

- Examples worked on board with illustration

Multiple Constraints

□ Now consider constraint: $g(\mathbf{w}) = [g_1(\mathbf{w}), \dots, g_M(\mathbf{w})]^T \leq 0$.

□ Lagrangian is:

$$L(\mathbf{w}, \boldsymbol{\lambda}) = f(\mathbf{w}) + \boldsymbol{\lambda}^T g(\mathbf{w}) = f(\mathbf{w}) + \sum_{m=1}^M \lambda_m g_m(\mathbf{w})$$

- Weighted sum of all M constraints
- $\boldsymbol{\lambda}$ is called the dual vector

□ KKT conditions extend to:

- $\hat{\mathbf{w}}$ minimizes the Lagrangian: $\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} L(\mathbf{w}, \hat{\boldsymbol{\lambda}})$
- For each $m = 1, \dots, M$
 - $g_m(\hat{\mathbf{w}}) = 0$ and $\hat{\lambda}_m \geq 0$ [active constraint]
 - $g_m(\hat{\mathbf{w}}) < 0$ and $\hat{\lambda}_m = 0$ [inactive constraint]

SVM Constrained Optimization

□ Recall: SVM constrained optimization

$$\min J_1(\mathbf{w}, b, \epsilon), \quad J_1(\mathbf{w}, b, \epsilon) = C \sum_{i=1}^N \epsilon_i + \frac{1}{2} \|\mathbf{w}\|^2$$

- Constraints: $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \epsilon_i$ and $\epsilon_i \geq 0$ for all $i = 1, \dots, N$

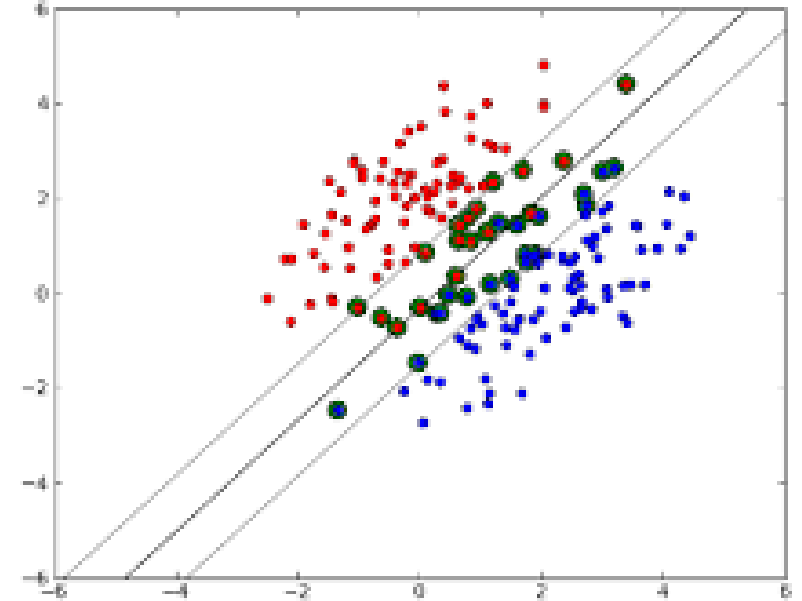
□ After applying KKT conditions and some algebra [beyond this class], solution is

- Optimal weight vector: $\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$ linear combination of instances
- Dual parameters α_i minimize

$$\sum_{i=1}^N \alpha_i + \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \quad \text{s.t.} \quad 0 \leq \alpha_i \leq C$$

Support Vectors

- Classifier weight is: $\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$
- Can show that $\alpha_i > 0$ only when \mathbf{x}_i is a support vector
 - On boundary or violating constraint
 - Otherwise $\alpha_i = 0$



What you should know

- ❑ Interpret weights in linear classification of images (logistic regression): Match filters
- ❑ Understand the margin in linear classification and maximum margin classifier
- ❑ SVM classifier: Allow violation of margin by introducing slack variables (More robust than linear classifier)
- ❑ Solve constrained optimization using the Lagrangian.
 - Understand KKT conditions for a single constraint
- ❑ Extend to nonlinear classifier by feature transformation: SVM with nonlinear kernels
- ❑ Select SVM parameters from cross-validation

Unit 8

K-NN and Kernel Methods

EE-UY 4563/EL-GY 9143: INTRODUCTION TO MACHINE LEARNING

Learning Objectives

- ❑ Mathematically describe K-NN and kernel methods
 - For regression and classification
- ❑ Visualize the regression and classification output
- ❑ Select kernel hyper-parameters from cross validation
- ❑ Compute the bias and variance error for estimators

Relation to SVM Lecture

- ❑ Unit 8 also has a lecture on SVM
- ❑ Beginning 2025, we are not teaching SVM
- ❑ Spending more time on Kernel methods and decision trees
- ❑ SVMs are not widely-used
- ❑ But look at SVM if you are interested

Outline



K -Nearest Neighbors

- ☐ Kernel Smoothing
- ☐ Kernel Regression
- ☐ Example with Climate Data



Nearest Neighbor Regression

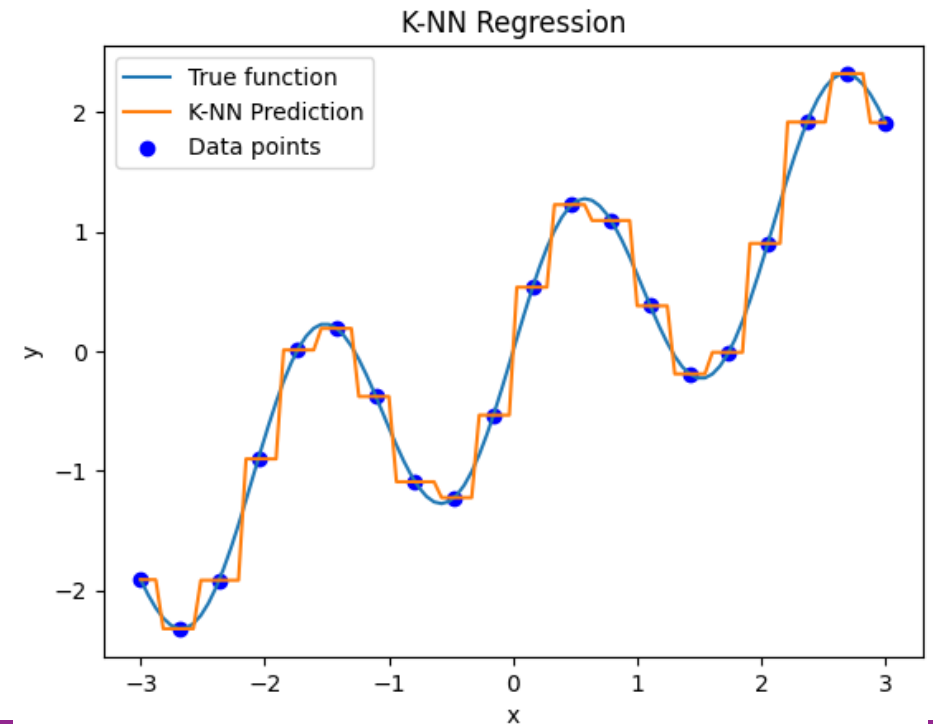
❑ **Problem:** Given data (x_i, y_i) make prediction for new x_0

❑ **Idea:** Similar values of x should have similar y

Nearest neighbor estimation

❑ Find x_i “closest” to x_0

❑ Select $\hat{y}_0 = y_i$



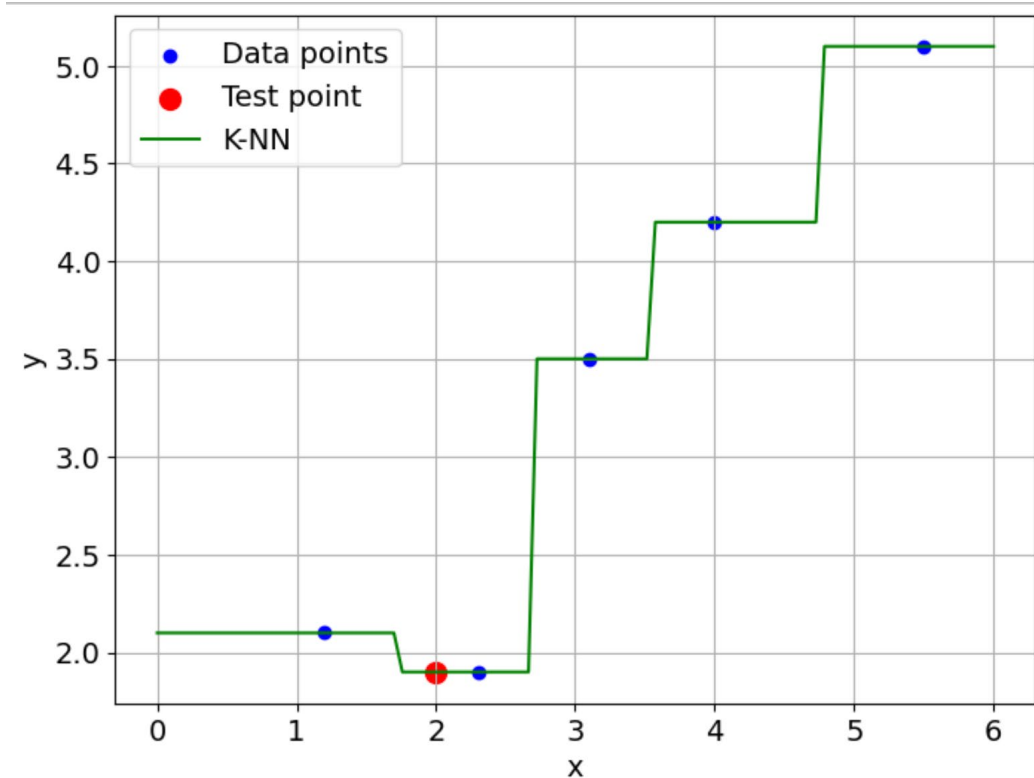
Example

Closest to
test point
 $x_0 = 2$

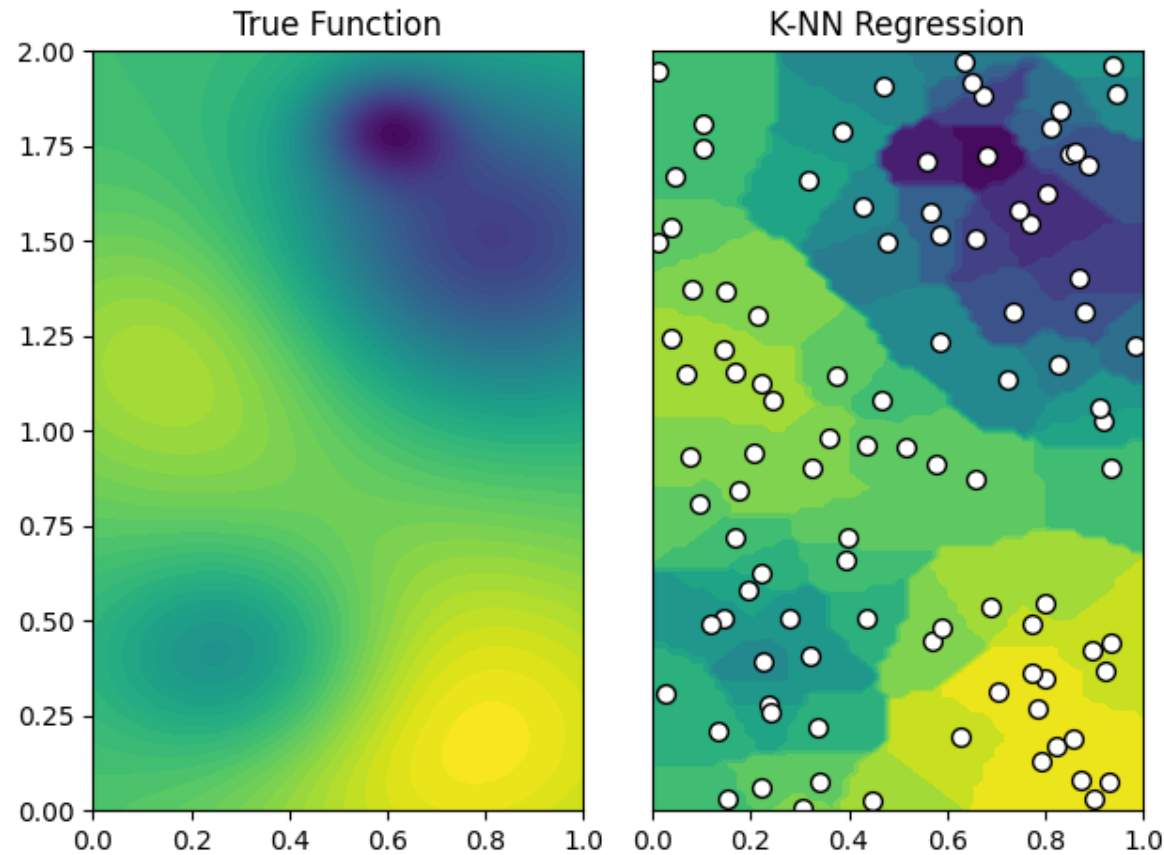
Data

i	x_i	y_i
1	1.2	2.1
2	2.3	1.9
3	3.1	3.5
4	4.0	4.2
5	5.5	5.1

Prediction
 $\hat{y}_0 = 1.9$



Nearest Neighbor in 2D



- Gives a “tessellation” of space
- Voronoi regions
 - Each region mapped to a point

K-Nearest Neighbor

❑ Problem with nearest neighbor:

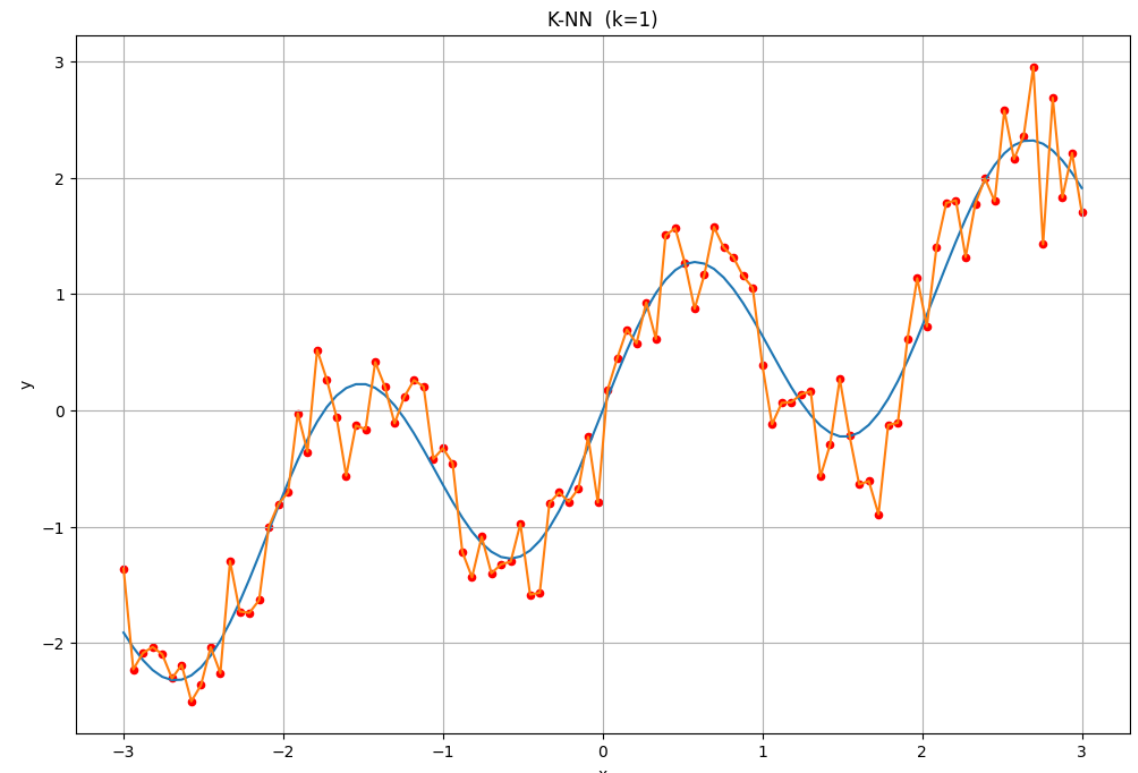
- Prediction sensitive to noise
- Overfits data

❑ K-Nearest neighbor:

- Take average of K closest points

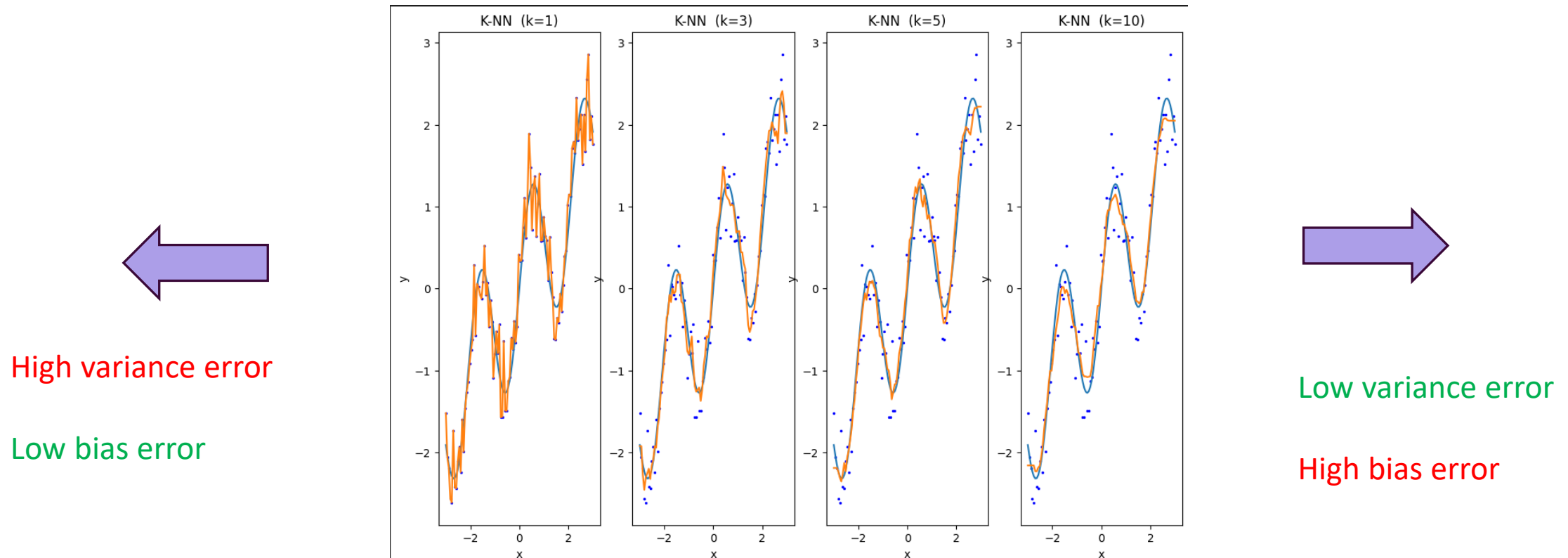
❑ Increasing K

- Reduces noise effect (lower variance error)
- Smooths (increase bias error)

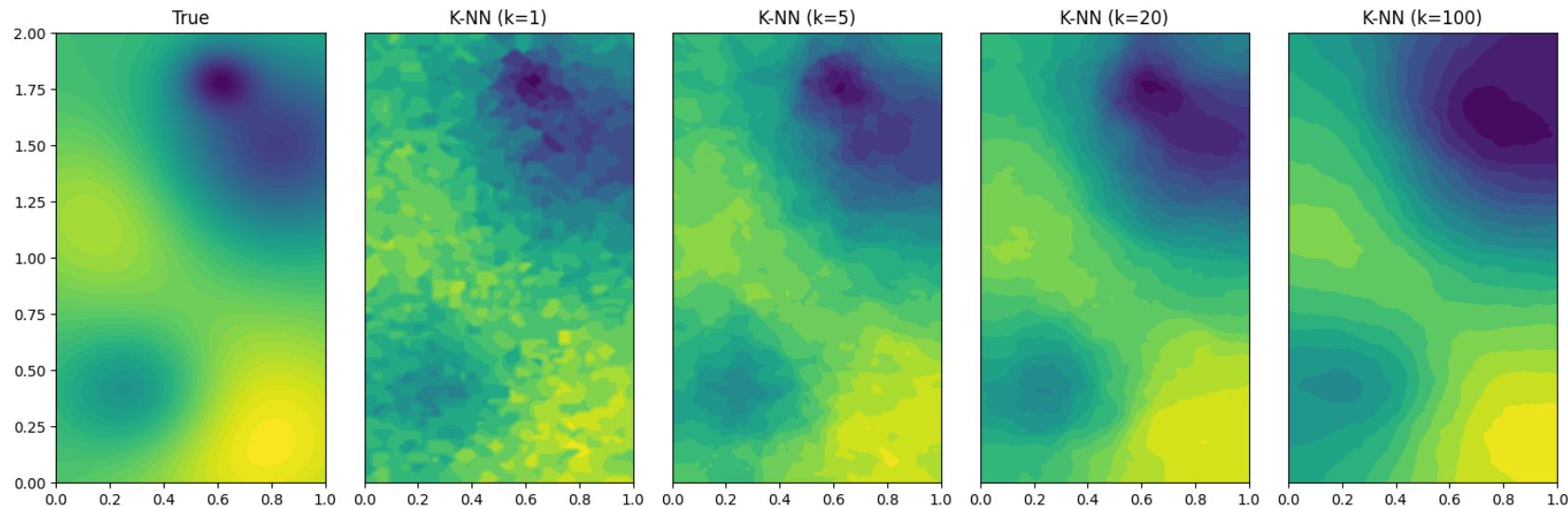


Example of overfitting data

1D Example



2D Example



High variance error
Low bias error

Low variance error
High bias error

Simulation details

- 1000 Data points
- $y_i = f(x_i) + w_i$
- $E|w_i|^2 = 0.3^2$

Error Analysis for K-NN

□ Suppose that $y_i = f(x_i) + w_i$,

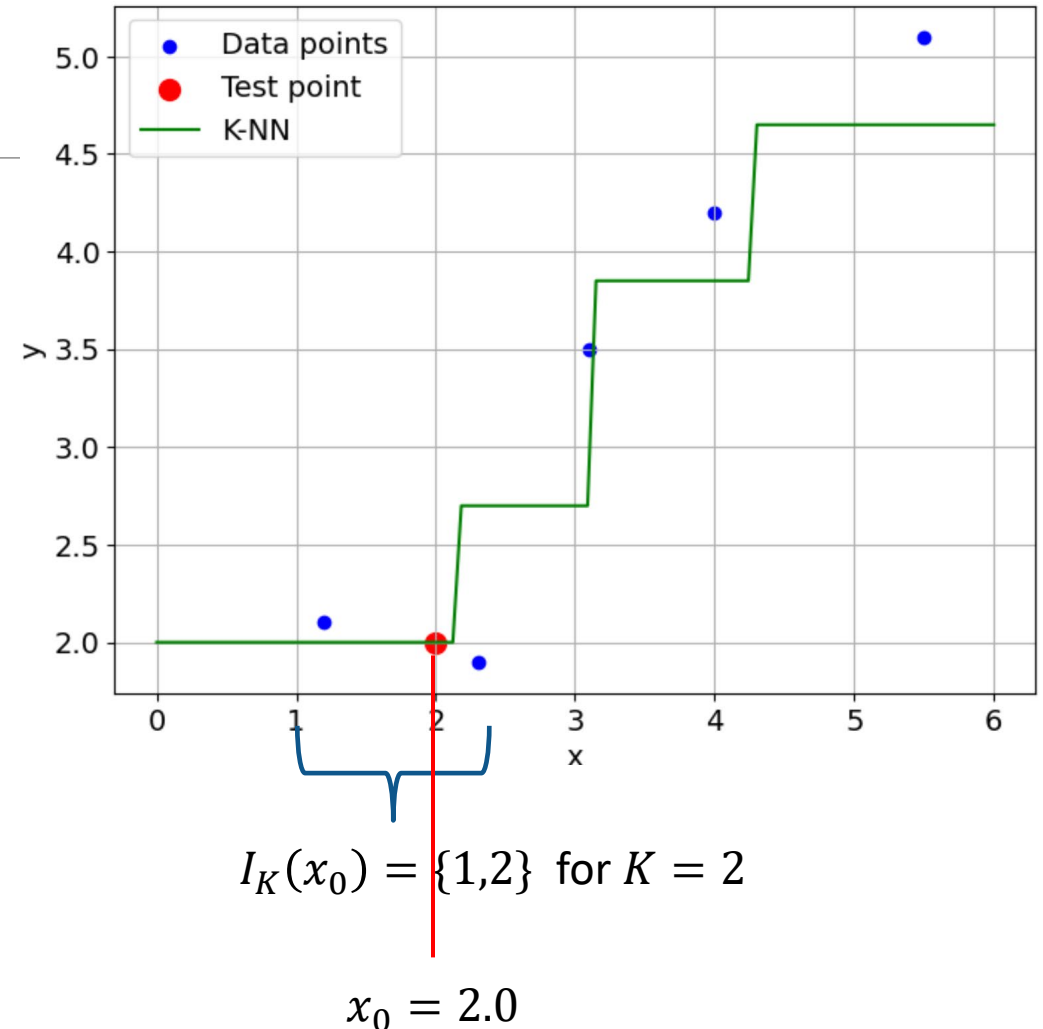
- w_i = measurement noise
- $E(w_i) = 0$, $E(w_i^2) = \sigma^2$

□ Given new test point x_0

□ Prediction is:

- $\hat{y}_0 = \frac{1}{K} \sum_{i \in I_K(x_0)} y_i$,
- $I_K(x_0) = K$ indices closest to x_0

□ What is the bias and variance error?



Bias and Variance

□ Prediction: $\hat{y}_0 = \frac{1}{K} \sum_{i \in I_K(x_0)} y_i$, $y_i = f(x_i) + w_i$

□ $E(\hat{y}_0) = \frac{1}{K} \sum_{i \in I_K(x_0)} f(x_i)$

□ Bias error: $Bias(x_0) = f(x_0) - \frac{1}{K} \sum_{i \in I_K(x_0)} f(x_i)$

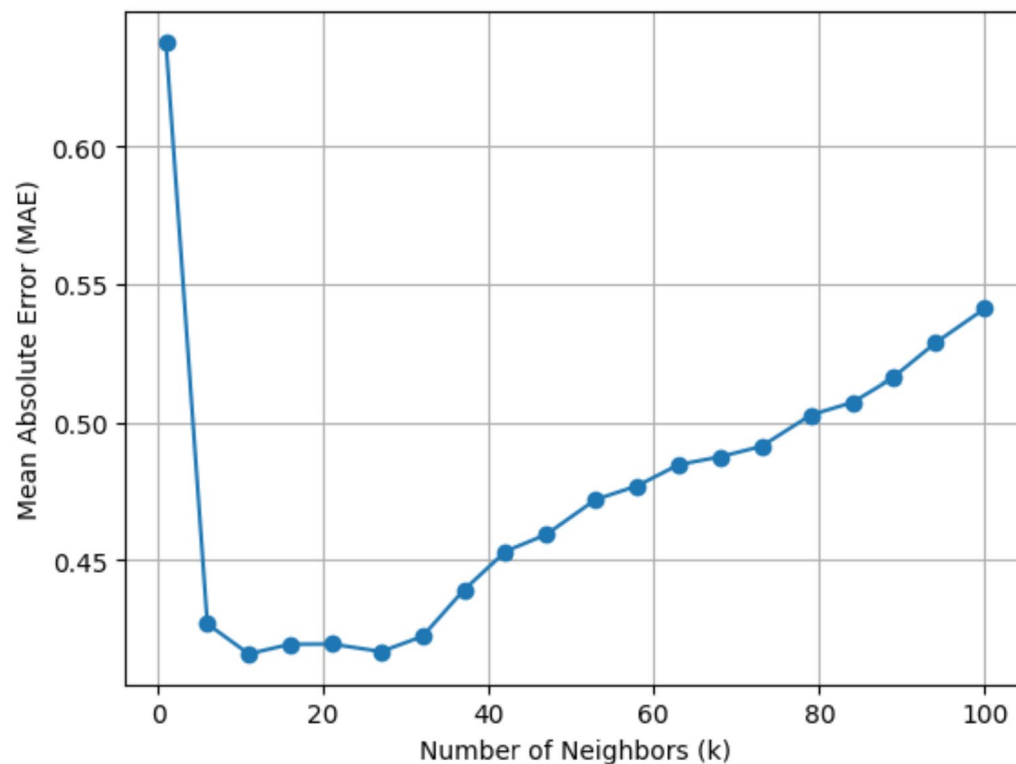
- Bias will be large when $f(x)$ varies over neighbors
- Increases with K

□ Variance error: $Var(x_0) = E[\hat{y}_0 - E(\hat{y}_0)]^2 = E\left[\frac{1}{K} \sum_{i \in I_K(x_0)} w_i\right]^2 = \frac{\sigma^2}{K}$

- Decreases with K

Using Cross-Validation

□ As usual find optimal K with cross-validation



```
mae = []
for i, k in enumerate(ktest):
    model = KNeighborsRegressor(n_neighbors=k)
    model.fit(xtrain, ytrain)
    yhat = model.predict(xtest)
    err = np.mean(np.abs(yhat - ytest))
    mae.append(err)
    print(f'k={k}, MAE={err:.3f}')
```

Parametric vs. Non-Parametric

Parametric methods

- ☐ Ex: Linear regression
- ☐ Assume function form
 - Model has parameters
- ☐ Fit parameters
- ☐ Use model directly

Non-Parametric methods

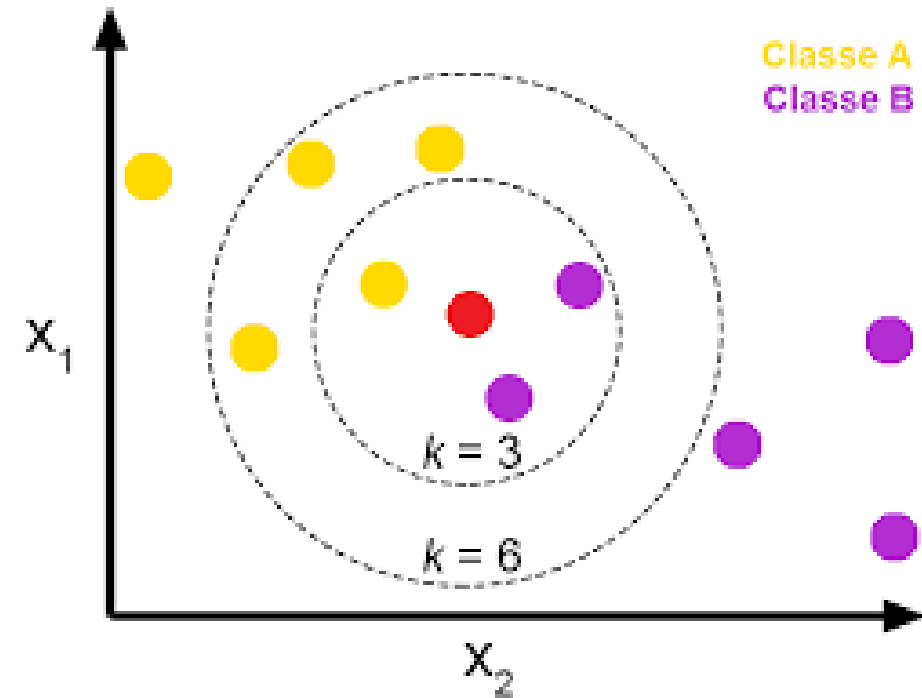
- ☐ Ex: K-NN
- ☐ No functional assumption
 - No parameters
 - But assumes smoothness
- ☐ No fitting
- ☐ Use data directly

Computational Cost Comparison

Task	Parametric Method (e.g. linear regression)	Non-Parametric (e.g. K-NN)
Model Fitting	Can be hard (e.g. gradient descent)	None, except hyper-parameters
Inference (after fitting)	Generally easier	Hard Must search through all data

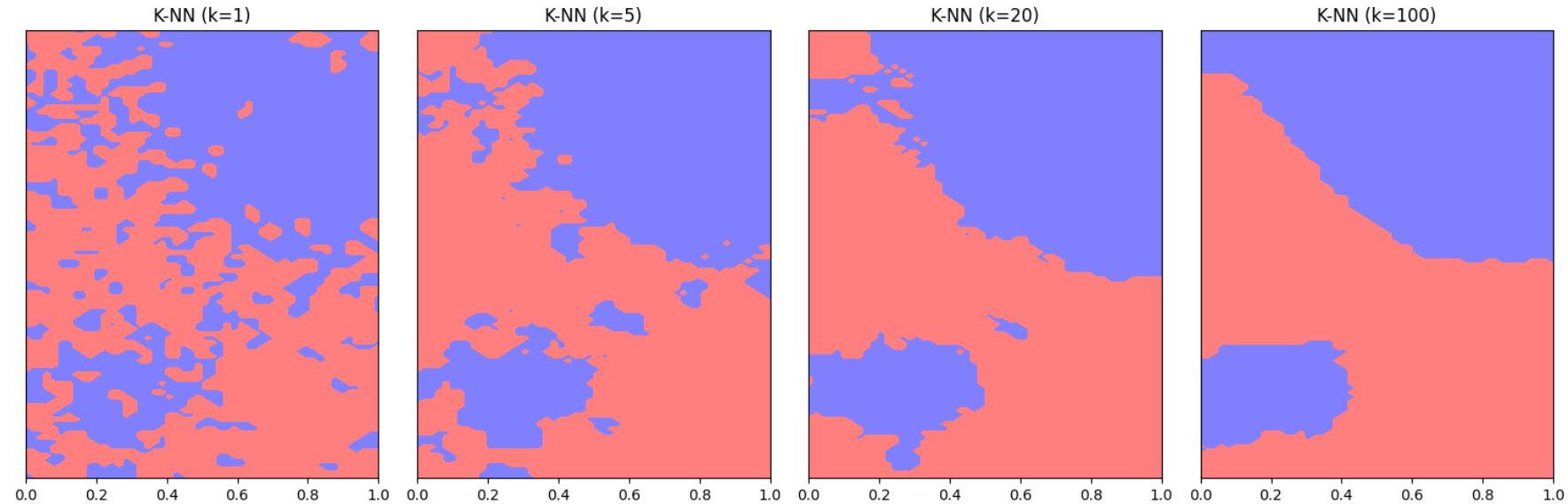
K-NN Classification

- K-NN can also be used for classification
- Data (x_i, y_i) binary labels $y_i \in \{0,1\}$
- Test point x_0
- Find indices of K closest neighbors $I_K(x_0)$
- Predict probability: $p(x_0) = \frac{1}{K} \sum_{i \in I_K(x_0)} y_i$
 - Fraction of neighbors that are class 1
- Class estimate $\hat{y}_0 = \begin{cases} 1 & p(x_0) \geq t \\ 0 & p(x_0) < t \end{cases}$
 - t = threshold. Typically, $t = \frac{1}{2}$



K-NN 2D Classification

- 1000 data points
- Larger K :
 - Smoother decision regions



Outline

☐ K -Nearest Neighbors

 ☒ Kernel Smoothing

☐ Kernel Regression

☐ Example with Climate Data

The Kernel Function

❑ Problem with K-NN

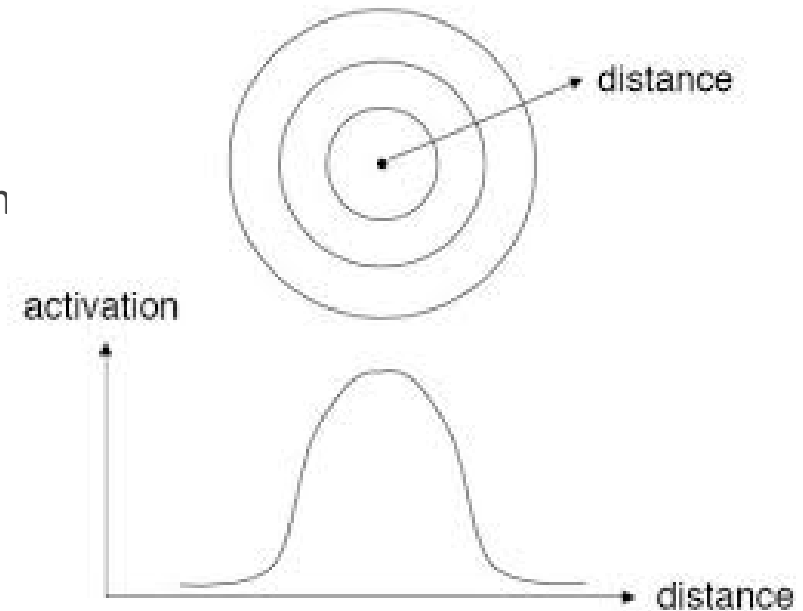
- Prediction is not smooth
- All K points weighted equally

❑ Kernel function:

- Function $K(x_i, x)$
- Measures “similarity” between new sample x and training

❑ Typical property

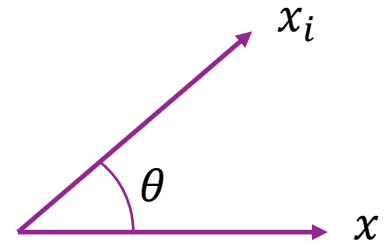
- x_i, x close $\Rightarrow K(x_i, x)$ maximum value
- x_i, x far $\Rightarrow K(x_i, x) \approx 0$



Common Kernels

Linear :

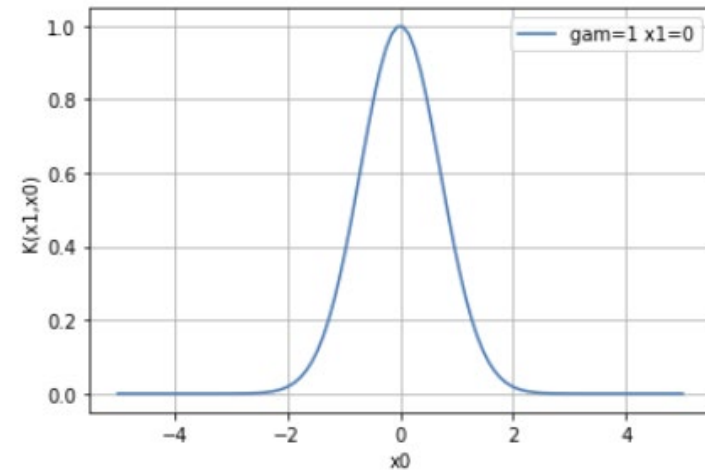
- $K(x_i, x) = x_i^T x = \|x_i\| \|x\| \cos \theta$
- Maximum when angle between vectors is small



Radial basis function:

$$K(x_i, x) = \exp[-\gamma \|x - x_i\|^2]$$

- $1/\gamma$ indicates width of kernel



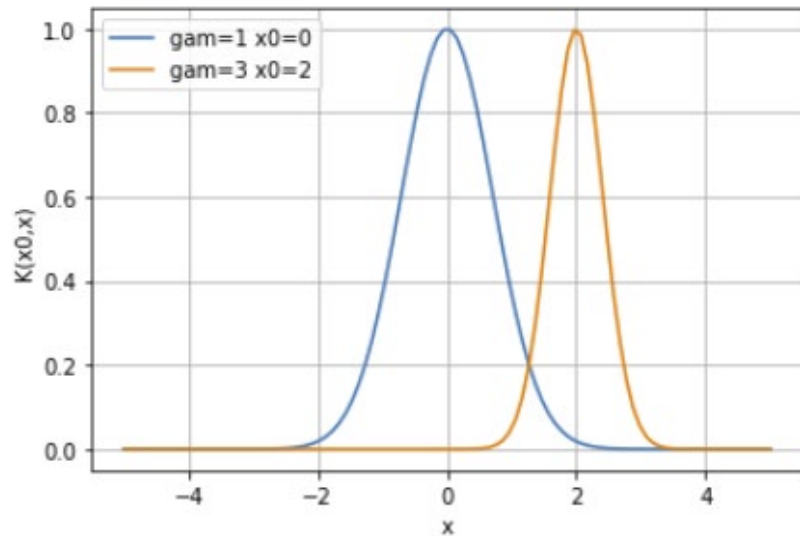
Polynomial kernel: $K(x_i, x) = |x_i^T x|^d$

- Typically $d=2$

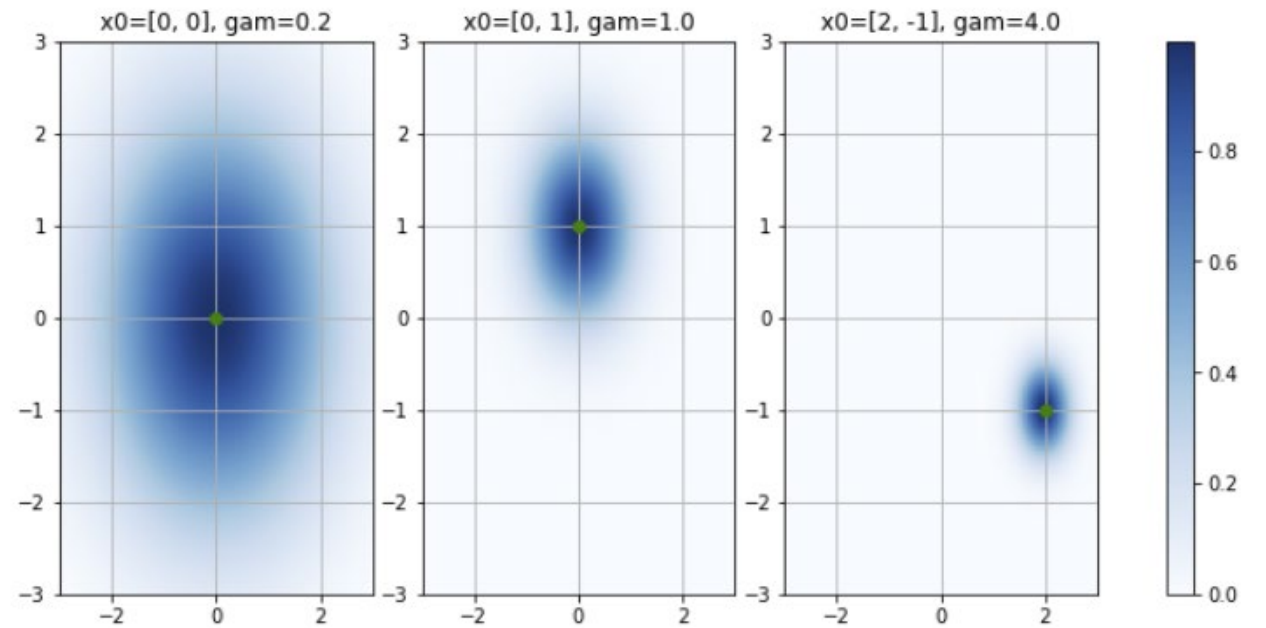
RBF Kernel Examples

□ RBF kernel: $K(x_0, x) = \exp[-\gamma \|x - x_0\|^2]$

- Peak value of 1 at $x = x_0$
- Width $\propto \frac{1}{\gamma}$



RBFs in 1D



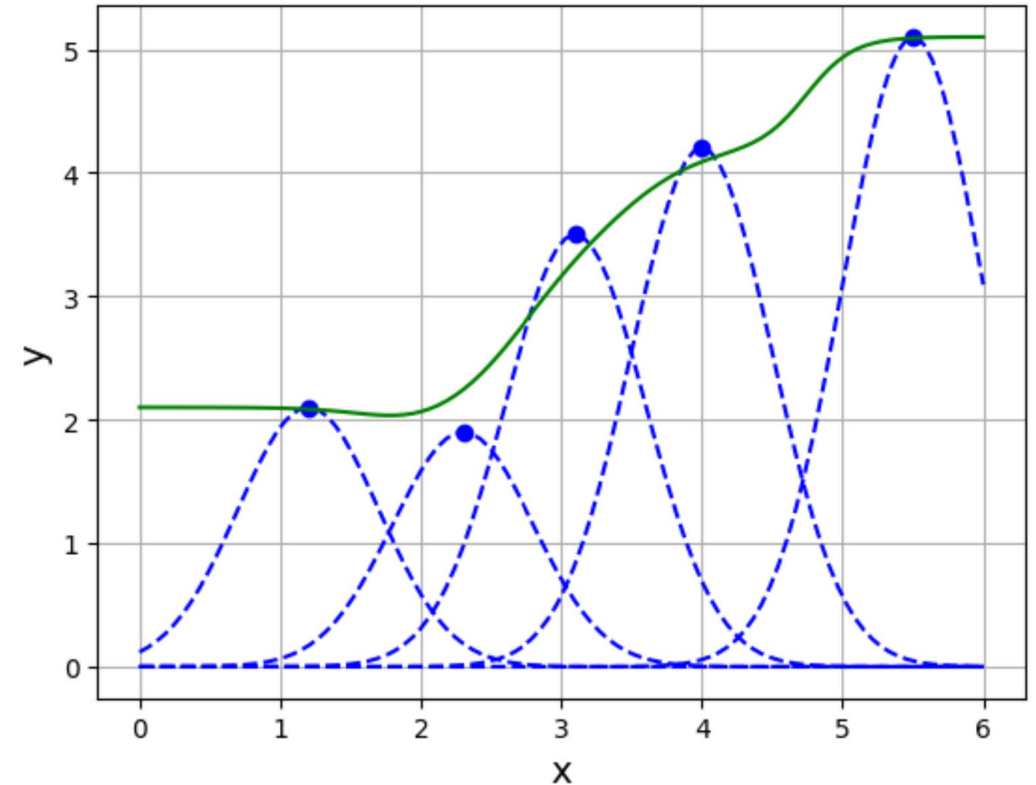
RBFs in 2D

Kernel Smoothing

- Given data $(x_i, y_i), i = 1, \dots, n$
- Target point x_0
- Compute distance to all data point $K(x_0, x_i)$
- Estimate at x_0

$$\hat{y}_0 = \frac{\sum_i y_i K(x_0, x_i)}{\sum_i K(x_0, x_i)}$$

- Weights each data point by $K(x_0, x_i)$
- Sometimes called **Nadaraya-Watson estimator**



Kernel Smoothing as Weighted Sum

□ Kernel smoother at x_0

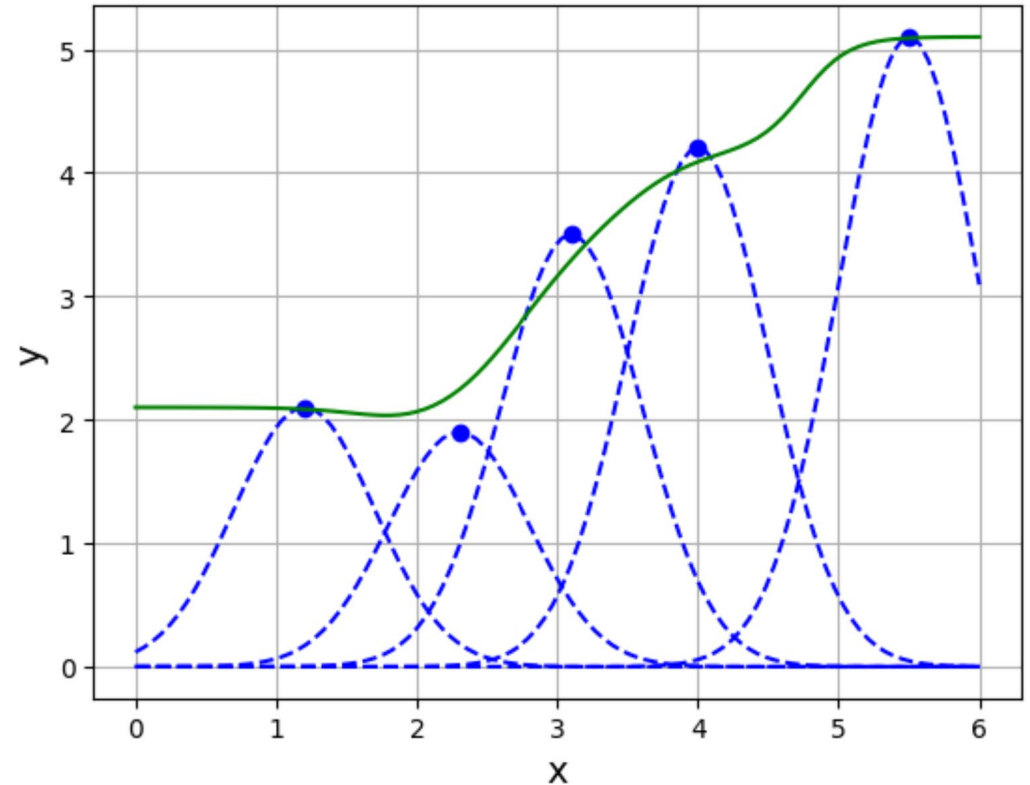
$$\hat{y}_0 = \frac{\sum_i y_i K(x_0, x_i)}{\sum_i K(x_0, x_i)} = \sum_i y_i \alpha_i$$

□ Weights: $\alpha_i = \frac{K(x_0, x_i)}{\sum_j K(x_0, x_j)}$

□ Note: $\sum_i \alpha_i = 1, \alpha_i \geq 0$

□ So, \hat{y}_0 is a weighted sum

□ Higher weight on points where $K(x_0, x_i)$



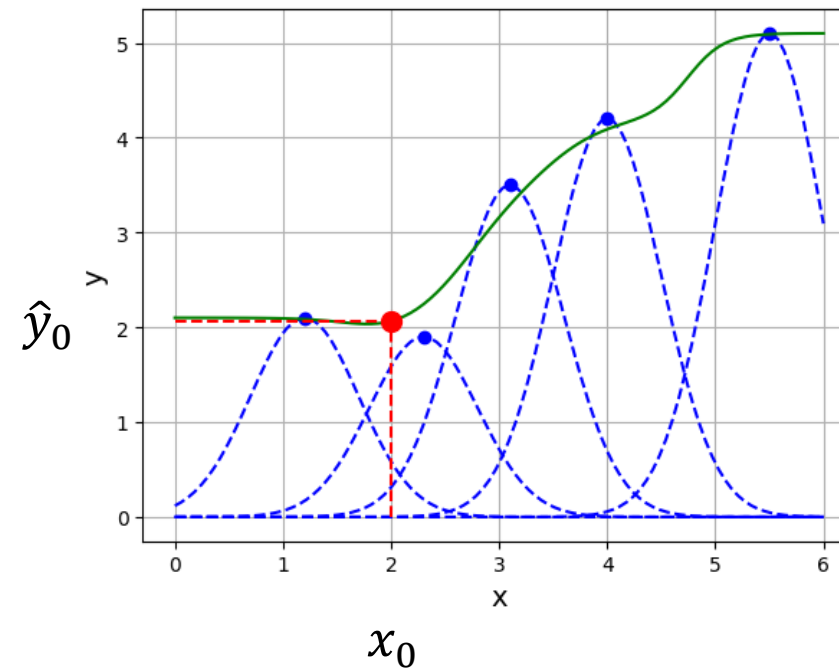
Example with RBF

Test point: $x_0 = 2$

RBF kernel: $\gamma = \frac{1}{\sigma^2}, \sigma = 0.5$

Estimate: $\hat{y}_0 = \sum_{i=1}^n \alpha_i x_i \approx 2.06$

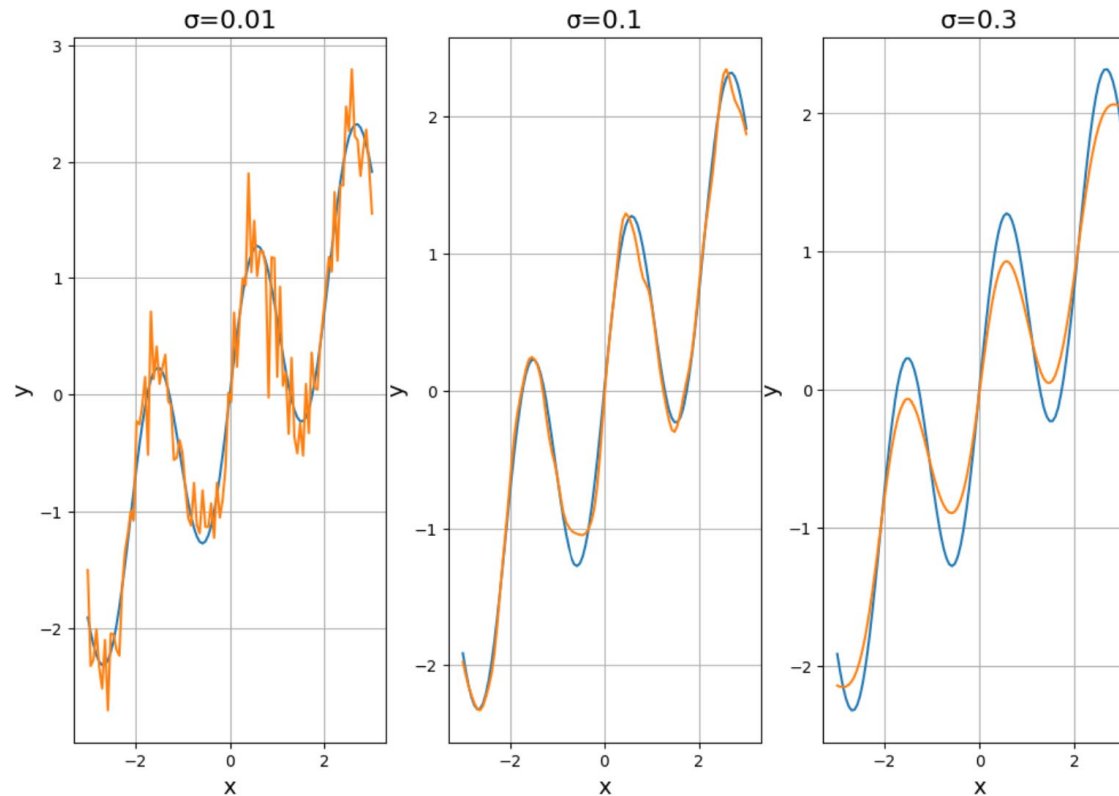
i	x_i	y_i	$K(x_0, x_i)$	α_i
1	1.2	2.1	0.27	0.23
2	2.3	1.9	0.84	0.69
3	3.1	3.5	0.089	0.074
4	4.0	4.2	$3.3(10)^{-4}$	$2.8(10)^{-4}$
5	5.5	5.1	$2.3(10)^{-11}$	$1.9(10)^{-11}$



1D Example

High variance error

Low bias error



Low variance error

High bias error

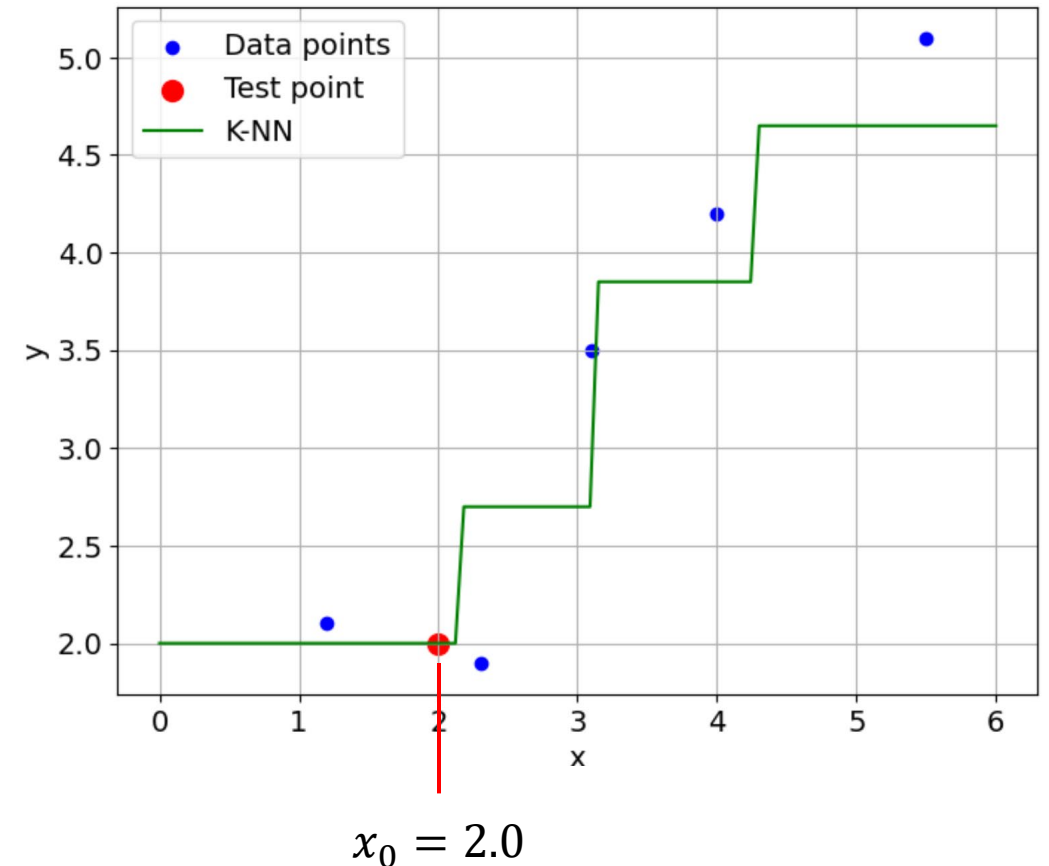
Kernel Smoother in Python

- ❑ No built-in function in sklearn
- ❑ But you can write it manually

```
def kernel_smoother(xdata, ydata, x, sig, tol=1e-8):  
    Dsq = np.sum((xdata[None, :, :] - x[:, None, :])**2, axis=2)  
    K = np.exp(-0.5 * Dsq / (sig**2))  
    Ksum = K.sum(axis=1, keepdims=True)  
    W = K / (Ksum + tol)  
    yhat = (W * ydata[None, :]).sum(axis=1)  
    return yhat
```

Error Analysis for Kernel

- Suppose that $y_i = f(x_i) + w_i$,
 - w_i = measurement noise
 - $E(w_i) = 0$, $E(w_i^2) = \sigma^2$
- Given new test point x_0
- Prediction is:
 - $\hat{y}_0 = \sum_i \alpha_i y_i$,
- What is the bias and variance error?



Bias and Variance

- Prediction: $\hat{y}_0 = \sum_i \alpha_i(x_0)y_i$, $y_i = f(x_i) + w_i$
- $E(\hat{y}_0) = \sum \alpha_i(x_0)f(x_i)$
- Bias error: $Bias(x_0) = f(x_0) - \sum \alpha_i(x_0)f(x_i)$
 - Bias will be large when $f(x)$ varies over kernel bandwidth
 - Increases with kernel radius
- Variance error: $Var(x_0) = E[\hat{y}_0 - E(\hat{y}_0)]^2 = E[\sum \alpha_i w_i]^2 = \sigma^2 \sum \alpha_i^2$

Variance Error for a Uniform Kernel

□ Example: Uniform kernel

$$K(x, x') = \begin{cases} 1 & \text{if } |x - x'| \leq \sigma \\ 0 & \text{else} \end{cases}$$

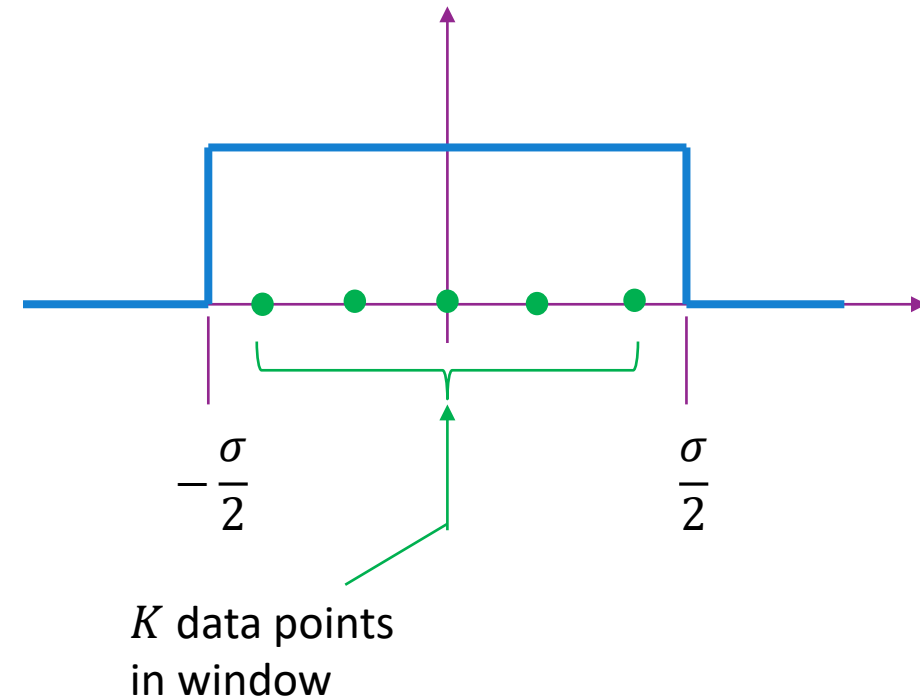
□ Data (x_i, y_i) and test point x_0

□ Suppose that M data points $|x_i - x_0| \leq \sigma$

□ Then: $\alpha_i(x_0) = \frac{K(x_0, x_i)}{\sum K(x_0, x_i)} = \frac{1}{M}$

□ Variance is $Var(x_0) = \sigma^2 \sum \alpha_i^2 = \frac{\sigma^2 M}{M^2} = \frac{\sigma^2}{M}$

□ Variance decreases with window size



Kernel Classifier

□ Kernel can also be used for classification

□ Given:

- Training data (x_i, y_i) with binary labels $y_i = \pm 1$
- Kernel $K(x_i, x)$

□ To classify a new point x :

- Decision function: $z = \sum_{i=1}^n y_i K(x_i, x)$
- Classify: $\hat{y} = \text{sign}(z)$

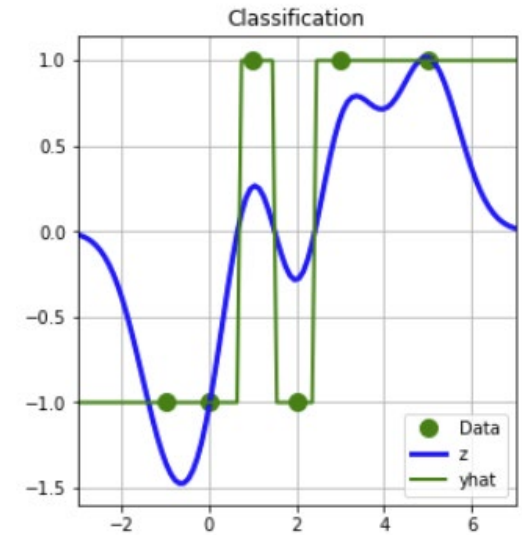
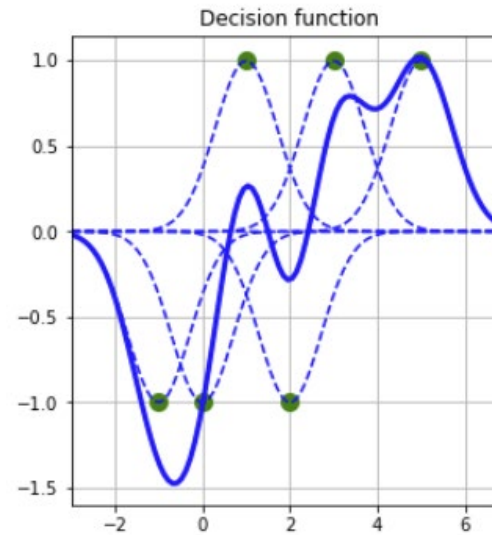
□ Idea:

- z is large positive when x is close to samples x_i with $y_i = 1$
- z is large negative when x is close to samples x_i with $y_i = -1$

Example in 1D

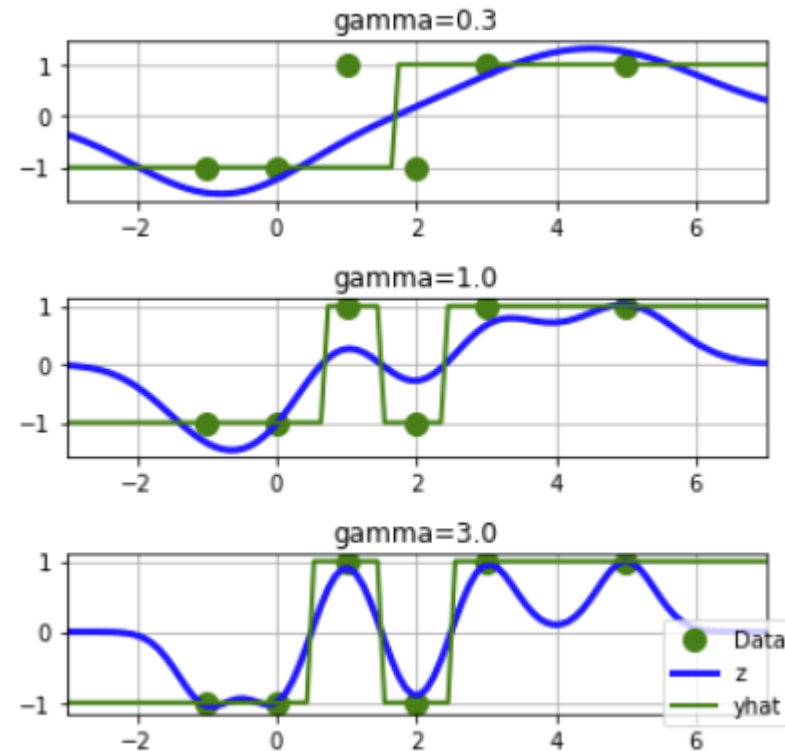
- Example data with 6 points (x_i, y_i)
 - RBF kernel: $K(x_i, x) = e^{-\gamma(x_i - x)^2}$, $\gamma = 1$
- Decision function:
 - $z = \sum_{i=1}^n y_i K(x_i, x)$
 - Sum of bell curves
 - Positive when near positive samples
 - Negative when near negative samples
- Classification:
 - $\hat{y} = \text{sign}(z)$

i	1	2	3	4	5	6
x_i	-1	0	1	2	3	5
y_i	-1	-1	1	-1	1	1



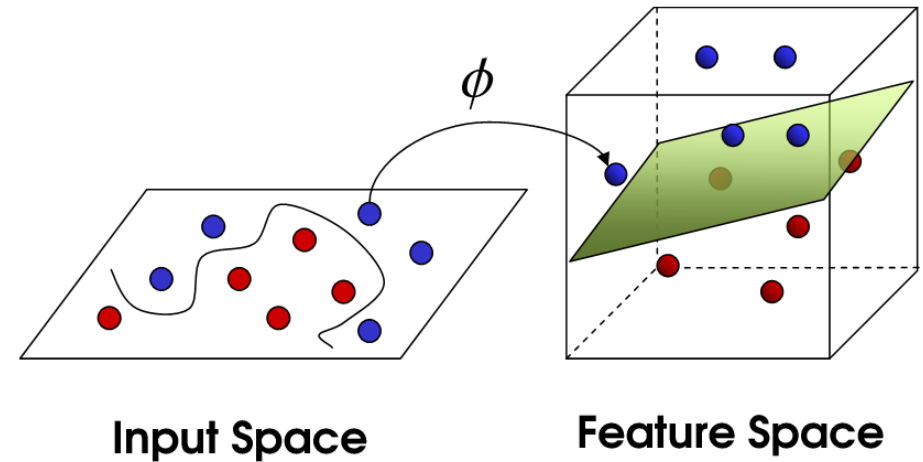
Effect of Gamma

- ❑ Same data as before
- ❑ RBF kernel: $K(x_i, x) = e^{-\gamma(x_i - x)^2}$
- ❑ As γ increases:
 - Decision function $z \approx y_i$ when $x = x_i$
 - Classifier fits training data better
 - Classification region more complex
- ❑ As a classifier, higher γ results in:
 - Lower bias error
 - But, higher variance error



Transform Linear Models

- Recall **transforms**
- Data (\mathbf{x}_i, y_i)
- Transform data $z_{ij} = \phi_j(\mathbf{x}_i), j = 1, \dots, p$
 - Each $\phi_j(\mathbf{x})$ is a “basis” function
- Linear model: $\hat{y} = \sum_{j=1}^p w_j \phi_j(\mathbf{x})$
- Weights w_j found by linear regression
- Maps original data \mathbf{x} to a higher dimensional space \mathbf{z}



Inner Product Kernel

- Consider a transform $\mathbf{z} = \phi(\mathbf{x})$
- Define kernel $K(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}')$
 - Inner product in the transformed space
- $K(\mathbf{x}, \mathbf{x}')$ has three key properties:
- Positive: $K(\mathbf{x}, \mathbf{x}) \geq 0$ since $K(\mathbf{x}, \mathbf{x}) = \|\phi(\mathbf{x})\|^2$
- Symmetry: $K(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \phi(\mathbf{x}') = \phi(\mathbf{x}')^T \phi(\mathbf{x}) = K(\mathbf{x}', \mathbf{x})$
- Positive semi-definite property: Given any α_i and \mathbf{x}_i $\sum_{i,j} \alpha_i \alpha_j K(\mathbf{x}_i, \mathbf{x}_j) \geq 0$
 - Why? $\sum_{i,j} \alpha_i \alpha_j K(\mathbf{x}_i, \mathbf{x}_j) = \sum_{i,j} \alpha_i \alpha_j \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j) = \left\| \sum_j \alpha_j \phi(\mathbf{x}_j) \right\|^2$

“Kernel Trick”

□ Consider transformed linear model: $\hat{y} = \hat{f}(\mathbf{x}) = \sum_{j=1}^p w_j \phi_j(\mathbf{x})$

□ Suppose \mathbf{w} solved with Ridge regression:

$$\mathbf{w} = \arg \min_{\mathbf{w}} \sum_i \left(y_i - \hat{f}(\mathbf{x}_i) \right)^2 + \lambda \|\mathbf{w}\|^2$$

□ **Theorem:** For any transformed linear model:

$$\hat{y} = \hat{f}(\mathbf{x}) = \sum_{i=1}^n \beta_i K(\mathbf{x}_i, \mathbf{x})$$

- Coefficients are solutions to $(\mathbf{K} + \lambda \mathbf{I})\boldsymbol{\beta} = \mathbf{y}$
- \mathbf{K} is the matrix with coefficients $K_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$

Implications of the Kernel Trick

- ❑ In any transformed model: $\hat{y} = \hat{f}(\mathbf{x}) = \sum_{i=1}^n \beta_i K(\mathbf{x}_i, \mathbf{x})$
- ❑ Similar to Kernel smoothing
 - But parameters β are fit from the data
- ❑ No need to transform data $\mathbf{x} \mapsto \phi(\mathbf{x})$
- ❑ Can solve for coefficients directly with Kernel $K(\mathbf{x}_i, \mathbf{x})$

Proof of Kernel Trick

- Weights are given by: $w = (A^T A + \lambda I)^{-1} A^T y$ with $A = \begin{bmatrix} \phi(x_1)^T \\ \vdots \\ \phi(x_n)^T \end{bmatrix}$
- From linear algebra $w = A^T (A A^T + \lambda I)^{-1} y$
- Let $\beta = (A A^T + \lambda I)^{-1} y$ so $w = A^T \beta$
- For any other data point x , $\hat{y} = \hat{f}(x) = w^T \phi(x) = \beta^T A^T \phi(x)$
- By definition of the kernel:
 - $[A A^T]_{ij} = \phi(x_i)^T \phi(x_j) = K(x_i, x_j)$ and $[A^T \phi(x)]_i = K(x_i, x)$
- Hence: $\hat{y} = \hat{f}(x) = w^T \phi(x) = \beta^T A^T \phi(x) = \sum \beta_i K(x_i, x)$
- Also: $(K + \lambda I)\beta = y$

Kernel Representation (Advanced)

- Say a kernel $K(\mathbf{x}, \mathbf{x}')$ is **symmetric** and **positive semidefinite** if:
 - Symmetry: $K(\mathbf{x}, \mathbf{x}') = K(\mathbf{x}', \mathbf{x})$
 - Given any α_i and \mathbf{x}_i $\sum_{i,j} \alpha_i \alpha_j K(\mathbf{x}_i, \mathbf{x}_j) \geq 0$
- Sometimes people say positive semidefinite to include symmetric
- **Theorem**: If kernel $K(\mathbf{x}, \mathbf{x}')$ is symmetric and positive semidefinite then there exists a mapping $\mathbf{x} \mapsto \phi(\mathbf{x})$ to a Hilbert space such that

$$K(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$$

- All kernels are inner products in some transformed space
- The space may be infinite dimensional

Outline

□ K -Nearest Neighbors

□ Kernel Smoothing

□ Kernel Regression

 Example with Climate Data

Temperature Modeling

- ❑ Problem: Interpolate temperature from limited measurements
- ❑ Data has natural spatial structure
- ❑ Kernel methods apply well
- ❑ Data from excellent [Berkeley Earth page](#)
- ❑ Several other sources as well

BERKELEY EARTH.

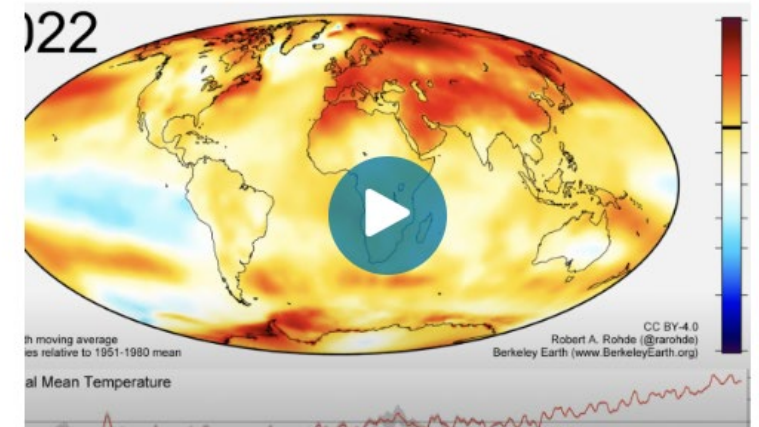
Data Visualization

Type

- ☐ Image
- ☐ Map
- ☐ Video

Subject

- ☐ Climate Science



2022 Global Temperature Anomaly Since 1850

Animation showing the change in global average temperature anomaly from 1850-2022 relative to the 1951- 1980 mean.

Loading the Data

- ❑ Get data from [data webpage](#)
- ❑ Can manually download file
- ❑ Or use downloader in jupyter notebook
- ❑ File is NetCDF (.nc) file
 - Used for multi-dimensional data

```
<xarray.Dataset> Size: 548MB
Dimensions:      (latitude: 180, longitude: 360, time: 2100, month_number: 12)
Coordinates:
  * longitude     (longitude) float32 1kB -179.5 -178.5 -177.5 ... 178.5 179.5
  * latitude      (latitude) float32 720B -89.5 -88.5 -87.5 ... 87.5 88.5 89.5
  * time          (time) float64 17kB 1.85e+03 1.85e+03 ... 2.025e+03 2.025e+03
Dimensions without coordinates: month_number
Data variables:
  land_mask      (latitude, longitude) float64 518kB ...
  temperature     (time, latitude, longitude) float32 544MB ...
  climatology     (month_number, latitude, longitude) float32 3MB ...
Attributes:
  Conventions:    Berkeley Earth Internal Convention (based on CF-1.5)
  title:          Native Format Berkeley Earth Surface Temperature A...
  history:        09-Jan-2025 20:35:17
  institution:    Berkeley Earth Surface Temperature Project
  land_source_history: 04-Jan-2025 19:11:11
  ocean_source_history: 06-Jan-2025 11:47:59
  comment:        This file contains Berkeley Earth surface temperat...
```

Global Gridded Data

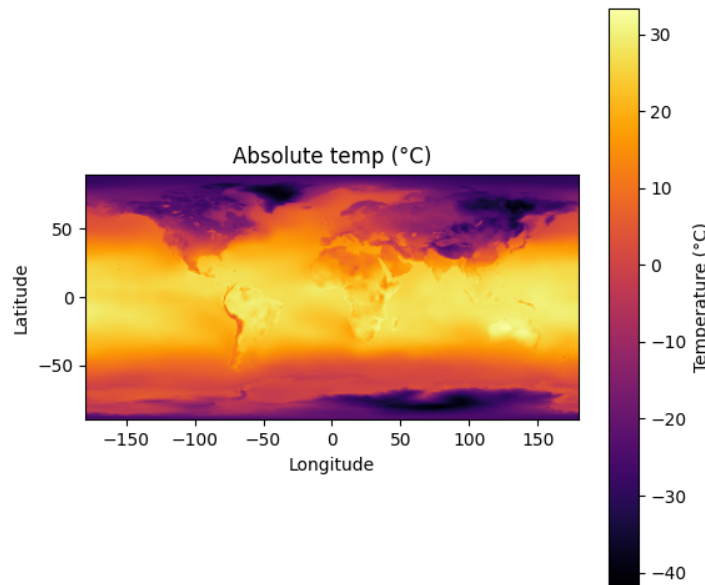
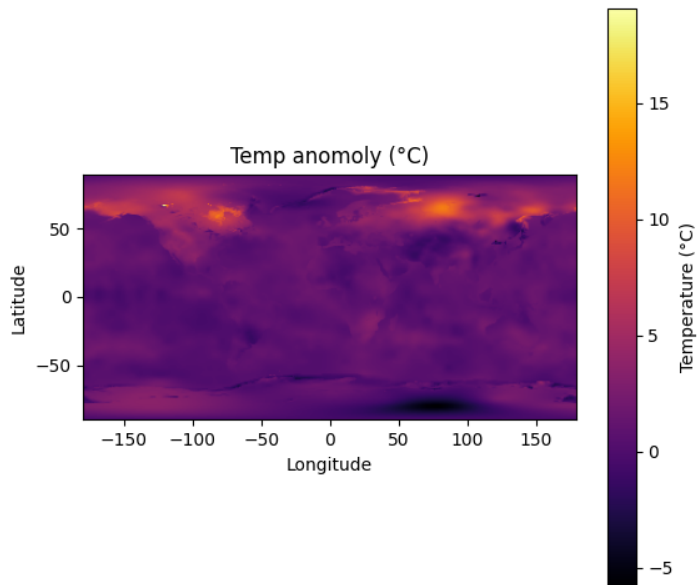
Datasets are also provided in a gridded NetCDF format. Two types of grids are provided, a grid based on a latitude-longitude grid. The equal area grid is the primary data format used in most of our analyses and may be less convenient for many users.

Datasets marked as “Experimental” below are products that are under development have not peer review give us feedback.

- Global Monthly Land + Ocean
 - Average Temperature with Air Temperatures at Sea Ice (Recommended; 1850 – Recent)
 - [Video of Temperature Field](#)
 - [Equal Area \(~100 MB\)](#)
 - [1° x 1° Latitude-Longitude Grid \(~400 MB\)](#)



Visualize the Temperature



- Berkeley data has climate
 - Over time and space
- We look at one “time slice”
- Temperature anomaly:
 - Difference from average monthly temp
- Temperature anomaly used for fitting
 - Removes known variation

Projections

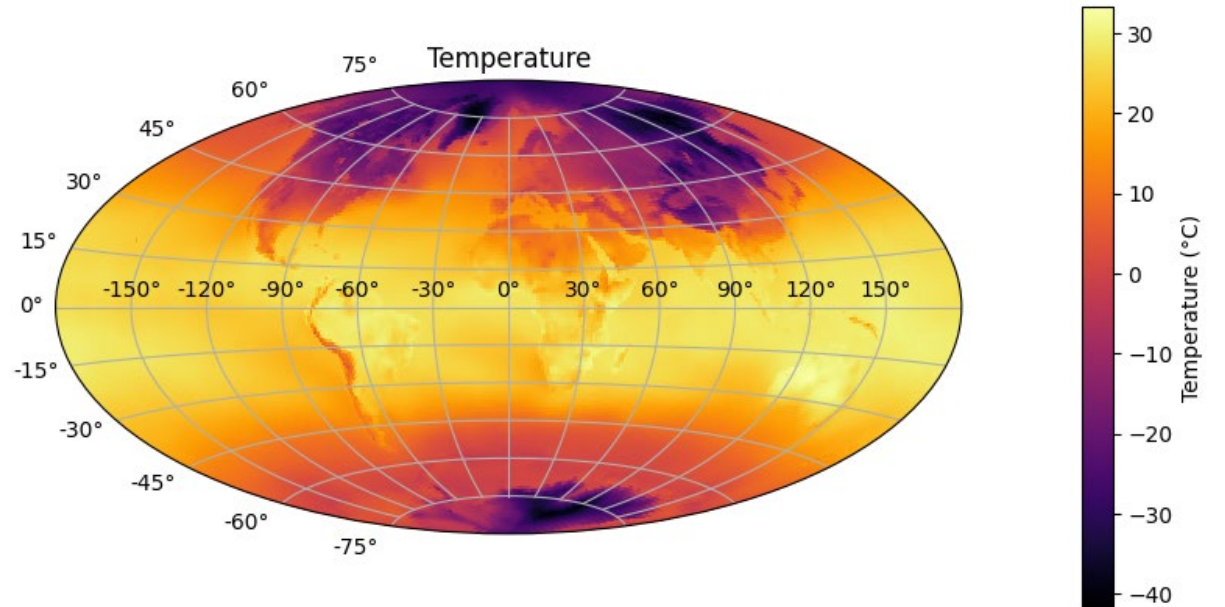
❑ You can also view via azimuthal map projection

❑ Example: Aitoff

```
# Convert coordinates to radians
lon_rad = np.deg2rad(lon) # shift to [-180, 180] then to radians
lat_rad = np.deg2rad(lat)

# Create meshgrid
lon_grid, lat_grid = np.meshgrid(lon_rad, lat_rad)

fig = plt.figure(figsize=(10, 5))
ax = fig.add_subplot(111, projection='aitoff')
im = ax.pcolormesh(lon_grid, lat_grid, temp_abs, cmap='inferno', shading='auto')
ax.grid(True)
plt.title('Temperature')
plt.colorbar(im, label='Temperature (°C)', pad=0.1)
```



Interpolation with K-NN

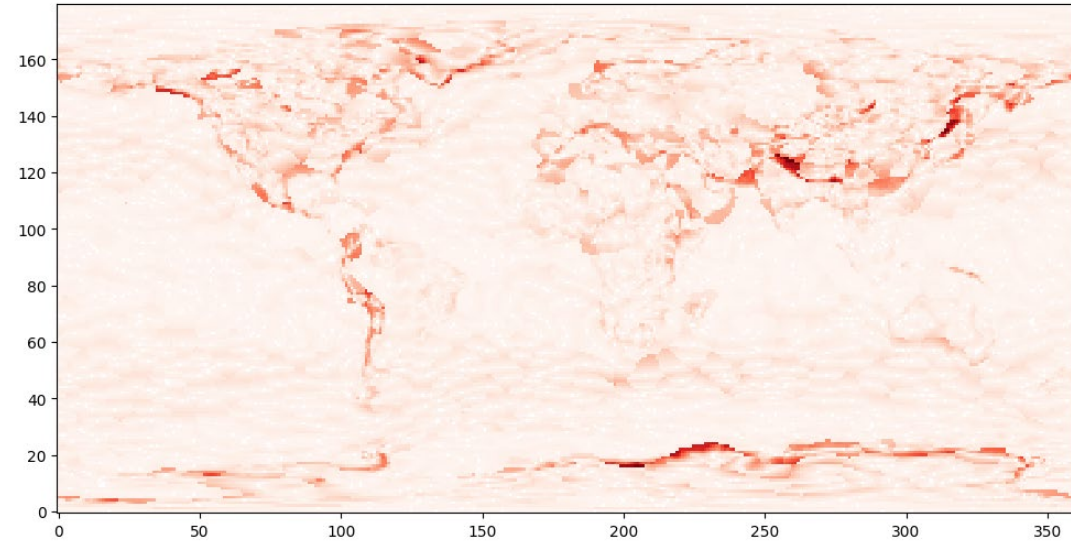
- ❑ Take 3% of the data for training
- ❑ Interpolate to remaining 90%
- ❑ Use (x, y, z) instead of lat-long for distance
 - Ensures distance is measured on sphere

```
# Convert the lat/long to Cartesian coordinates for distance calculation
# We use a unit sphere (R=1) for simplicity
R = 1
x = R * np.cos(np.deg2rad(lat_flat)) * np.cos(np.deg2rad(lon_flat))
y = R * np.cos(np.deg2rad(lat_flat)) * np.sin(np.deg2rad(lon_flat))
z = R * np.sin(np.deg2rad(lat_flat))

# Select a random subset of points for training
p = 0.03 # fraction of points to use for training
num_points = len(temp_flat)
num_train = int(p * num_points)
np.random.seed(0) # for reproducibility
train_indices = np.random.choice(num_points, num_train, replace=False)
test_indices = np.setdiff1d(np.arange(num_points), train_indices)
X_train = np.vstack((x[train_indices], y[train_indices], z[train_indices])).T
y_train = temp_flat[train_indices]
X_test = np.vstack((x[test_indices], y[test_indices], z[test_indices])).T
y_test = temp_flat[test_indices]
```

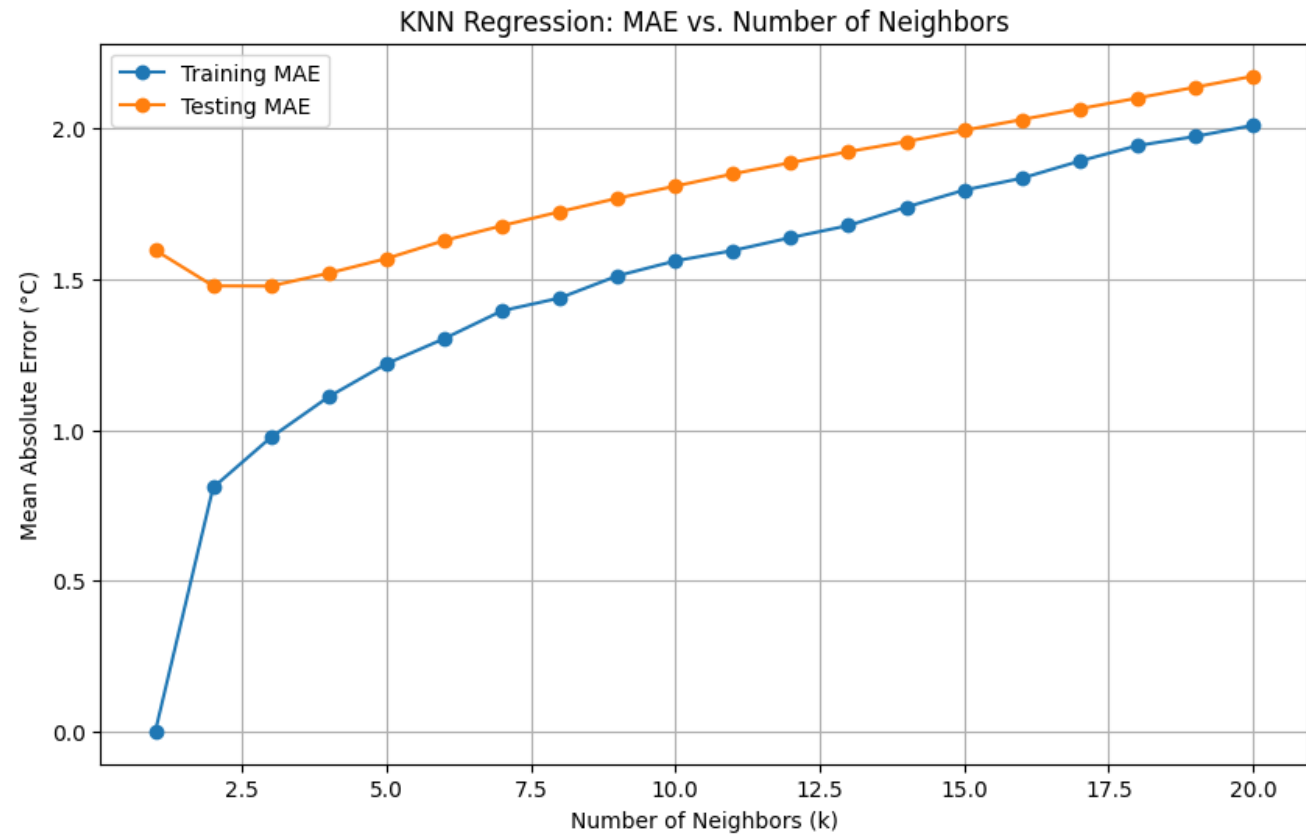

Interpolation with K-NN

- ❑ Error with $K = 1$
- ❑ MAE = 1.6°
- ❑ Error map to the right
 - Errors largest at mountains and coastlines
 - Locations of rapid variation

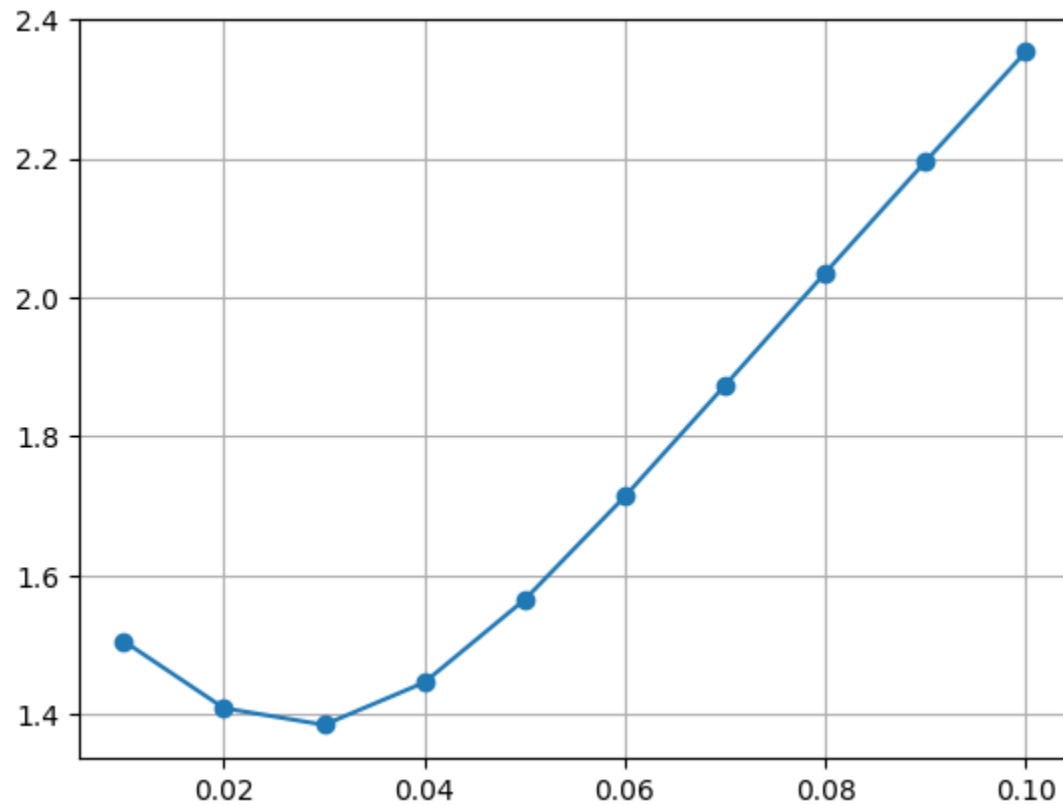


Selecting K

□ $K = 3$ is optimal



Kernel Smoother



□ Compute MAE vs. bandwidth

□ Performs slightly better than K-NN