

# Nearest Neighbors (continued), Perceptron

Sriram Sankararaman

The instructor gratefully acknowledges Fei Sha, Ameet Talwalkar, Eric Eaton, and Jessica Wu whose slides are heavily used, and the many others who made their course material freely available online.

# Announcements

- If you attempted mini-quiz and problem set 0, you should have received a PTE.
- If not, talk to one of the TAs at the end of lecture.

# Outline

- 1 Review of previous lecture
  - Nearest neighbor classifier
  - Some practical sides of NNC
  - Preprocessing data
- 2 Perceptron
- 3 What we have learned

# Nearest neighbor classification (NNC)

## Training

- Store the entire training set.

# Nearest neighbor classification (NNC)

## Testing

$$\mathbf{x}(1) = \mathbf{x}_{\text{nn}(\mathbf{x})}$$

where  $\text{nn}(\mathbf{x}) \in [\mathbf{N}] = \{1, 2, \dots, \mathbf{N}\}$ , i.e., the index to one of the training instances

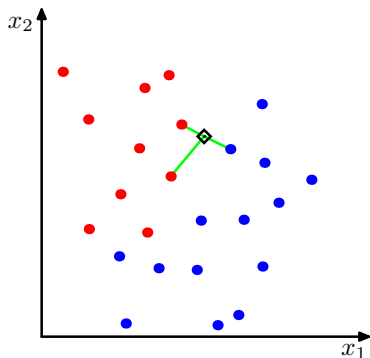
$$\text{nn}(\mathbf{x}) = \arg \min_{n \in [\mathbf{N}]} \|\mathbf{x} - \mathbf{x}_n\|_2^2 = \arg \min_{n \in [\mathbf{N}]} \sum_{d=1}^D (x_d - x_{nd})^2$$

## Classification rule

$$y = h(\mathbf{x}) = y_{\text{nn}(\mathbf{x})}$$

## Visual example

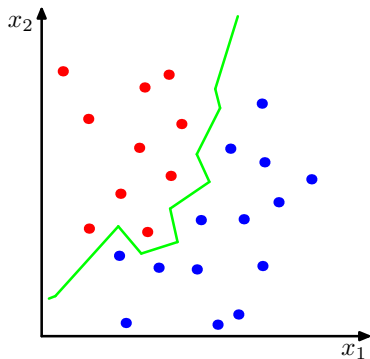
In this 2-dimensional example, the nearest point to  $x$  is a **red training instance**, thus,  $x$  will be labeled as **red**.



(a)

# Decision boundary

For every point in the space, determine its label using the NNC rule. Gives rise to a *decision boundary* that partitions the space into different regions.



(b)





# K-nearest neighbor (KNN) classification

## Increase the number of nearest neighbors to use?

- 1-nearest neighbor:  $nn_1(\mathbf{x}) = \arg \min_{n \in [N]} \|\mathbf{x} - \mathbf{x}_n\|_2^2$
- 2nd-nearest neighbor:  $nn_2(\mathbf{x}) = \arg \min_{n \in [N] - nn_1(\mathbf{x})} \|\mathbf{x} - \mathbf{x}_n\|_2^2$
- 3rd-nearest neighbor:  $nn_3(\mathbf{x}) = \arg \min_{n \in [N] - nn_1(\mathbf{x}) - nn_2(\mathbf{x})} \|\mathbf{x} - \mathbf{x}_n\|_2^2$

# K-nearest neighbor (KNN) classification

## Increase the number of nearest neighbors to use?

- 1-nearest neighbor:  $\text{nn}_1(\mathbf{x}) = \arg \min_{n \in [\mathbf{N}]} \|\mathbf{x} - \mathbf{x}_n\|_2^2$
- 2nd-nearest neighbor:  $\text{nn}_2(\mathbf{x}) = \arg \min_{n \in [\mathbf{N}] - \text{nn}_1(\mathbf{x})} \|\mathbf{x} - \mathbf{x}_n\|_2^2$
- 3rd-nearest neighbor:  $\text{nn}_3(\mathbf{x}) = \arg \min_{n \in [\mathbf{N}] - \text{nn}_1(\mathbf{x}) - \text{nn}_2(\mathbf{x})} \|\mathbf{x} - \mathbf{x}_n\|_2^2$

## The set of K-nearest neighbors

$$\text{knn}(\mathbf{x}) = \{\text{nn}_1(\mathbf{x}), \text{nn}_2(\mathbf{x}), \dots, \text{nn}_K(\mathbf{x})\}$$

Let  $\mathbf{x}(k) = \mathbf{x}_{\text{nn}_k(\mathbf{x})}$ , then

$$\|\mathbf{x} - \mathbf{x}(1)\|_2^2 \leq \|\mathbf{x} - \mathbf{x}(2)\|_2^2 \leq \dots \leq \|\mathbf{x} - \mathbf{x}(K)\|_2^2$$

# How to classify with $K$ neighbors?

# How to classify with $K$ neighbors?

## Classification rule

- Every neighbor votes: suppose  $y_n$  (the label) for  $\mathbf{x}_n$  is  $c$ , then
  - ▶ vote for  $c$  is 1
  - ▶ vote for  $c' \neq c$  is 0

We use the *indicator function*  $\mathbb{I}(y_n == c)$  to represent.

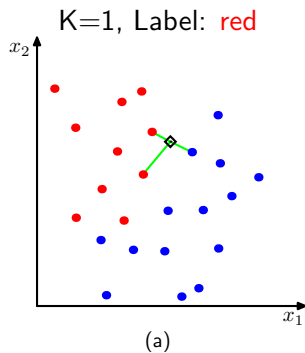
- Aggregate everyone's vote

$$v_c = \sum_{n \in \text{knn}(\mathbf{x})} \mathbb{I}(y_n == c), \quad \forall \quad c \in [\mathbf{C}]$$

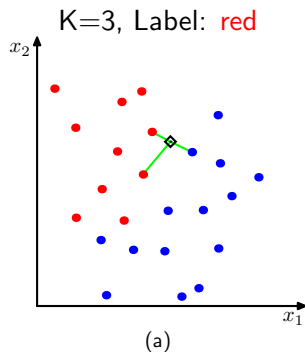
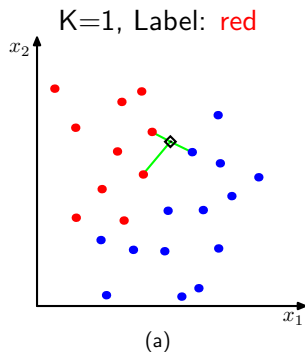
- Label with the majority

$$y = h(\mathbf{x}) = \arg \max_{c \in [\mathbf{C}]} v_c$$

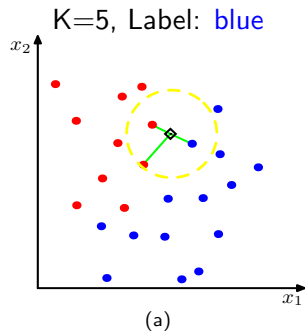
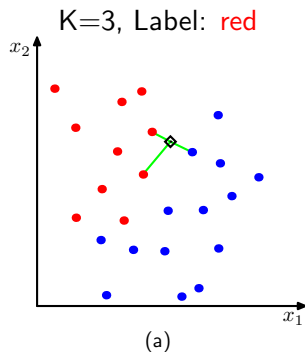
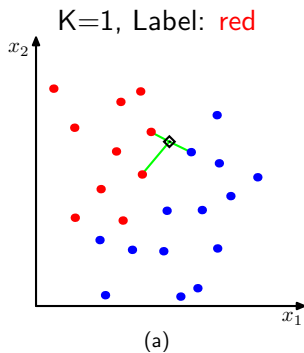
# Example



# Example



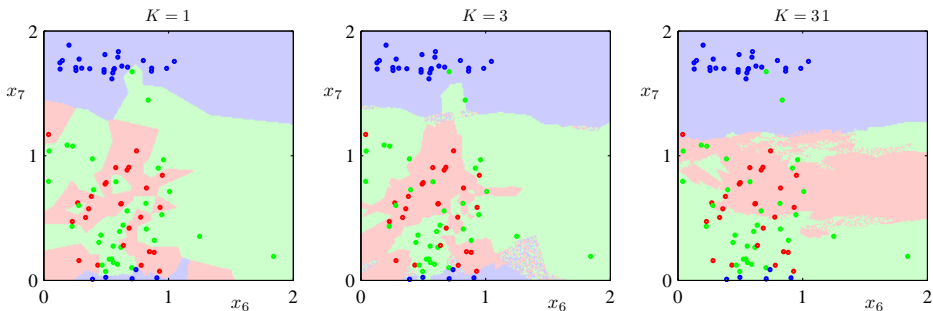
# Example



# Decision boundary as a function of $K$



# Decision boundary as a function of $K$



When  $K$  increases, the decision boundary becomes smooth.

# Mini-summary

## Advantages of NNC

- Computationally, simple and easy to implement – just computing the distance

## Disadvantages of NNC

- Computationally intensive for large-scale problems:  $O(ND)$  for labeling a data point
- We need to “carry” the training data around. Without it, we cannot do classification. This type of method is called *nonparametric*.
- Choosing the right distance measure and  $K$  can be involved.

# Hypeparameters in NNC

## Two practical issues about NNC

- Choosing  $K$ , i.e., the number of nearest neighbors (default is 1)
- Choosing the right distance measure (default is Euclidean distance), for example, from the following generalized distance measure

$$\|\mathbf{x} - \mathbf{x}_n\|_p = \left( \sum_d |x_d - x_{nd}|^p \right)^{1/p}$$

for  $p \geq 1$ .

*Those are not specified by the algorithm itself — resolving them requires empirical studies and are task/dataset-specific.*

# Tuning by using a validation dataset

## Training data (set)

- N samples/instances:  $\mathcal{D}^{\text{TRAIN}} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$
- They are used for learning  $h(\cdot)$

## Test (evaluation) data

- M samples/instances:  $\mathcal{D}^{\text{TEST}} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_M, y_M)\}$
- They are used for assessing how well  $h(\cdot)$  will do in predicting an unseen  $\mathbf{x} \notin \mathcal{D}^{\text{TRAIN}}$

## Development (or validation) data

- L samples/instances:  $\mathcal{D}^{\text{DEV}} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_L, y_L)\}$
- They are used to optimize hyperparameter(s).

Training data, validation and test data should *not* overlap!

# Recipe

- For each possible value of the hyperparameter (say  $K = 1, 3, \dots, 100$ )
  - ▶ Train a model using  $\mathcal{D}^{\text{TRAIN}}$
  - ▶ Evaluate the performance of the model on  $\mathcal{D}^{\text{DEV}}$
- Choose the model with the best performance on  $\mathcal{D}^{\text{DEV}}$
- Evaluate the model on  $\mathcal{D}^{\text{TEST}}$

# What if we do not have validation data?

- Use cross-validation.

# Yet, another practical issue with NNC

## **Assumes all features are equally important!**

- Distances depend on units of the features.

## Example: classify Iris with two features

### Training data

ID (n)	petal width ( $x_1$ )	sepal length ( $x_2$ )	category ( $y$ )
1	0.2	5.1	setosa
2	1.4	7.0	versicolor
3	2.5	6.7	virginica

### Flower with unknown category

petal width = 1.8 and sepal width = 6.4

Calculating distance =  $\sqrt{(x_1 - x_{n1})^2 + (x_2 - x_{n2})^2}$

ID	distance
1	1.75
2	0.72
3	0.76

Thus, the category is *versicolor* (the real category is *virginica*)



Change units of  $x_2$  from *cm* to *mm*

## Training data

ID (n)	petal width ( $x_1$ )	sepal length ( $x_2$ )	category ( $y$ )
1	0.2	51	setosa
2	1.4	70	versicolor
3	2.5	67	virginica

## Change units of $x_2$ from $cm$ to $mm$

### Training data

ID (n)	petal width ( $x_1$ )	sepal length ( $x_2$ )	category ( $y$ )
1	0.2	51	setosa
2	1.4	70	versicolor
3	2.5	67	virginica

### Flower with unknown category

petal width = 1.8 and sepal width = 64

Calculating distance =  $\sqrt{(x_1 - x_{n1})^2 + (x_2 - x_{n2})^2}$

ID	distance
1	13
2	6
3	3

Thus, the category is *virginica* (the real category is *virginica*)

# Preprocess data

## Normalize data to have zero mean and unit standard deviation in each dimension

- Compute the means and standard deviations in each feature

$$\bar{x}_d = \frac{1}{N} \sum_n x_{nd}, \quad s_d^2 = \frac{1}{N-1} \sum_n (x_{nd} - \bar{x}_d)^2$$

- Scale the feature accordingly

$$x_{nd} \leftarrow \frac{x_{nd} - \bar{x}_d}{s_d}$$

*Many other ways of normalizing data — you would need/want to try different ones and pick them using (cross)validation*

# Outline

- 1 Review of previous lecture
- 2 Perceptron
  - Setup for binary classification
  - Intuition
  - Algorithm
- 3 What we have learned

# Perceptron learning

## Special case: binary classification

- Instance (**feature vectors**):  $\mathbf{x} \in \mathbb{R}^D$
- **Label**:  $y \in \{-1, +1\}$
- **Model/Hypotheses**:  
$$H = \{h | h : \mathbb{X} \rightarrow \mathbb{Y}, h(\mathbf{x}) = \text{sign}(\sum_{d=1}^D w_d x_d + b)\}.$$
- Learning goal:  $\hat{y} = h(\mathbf{x})$

# Perceptron learning

## Special case: binary classification

- Instance (**feature vectors**):  $\mathbf{x} \in \mathbb{R}^D$
- **Label**:  $y \in \{-1, +1\}$
- **Model/Hypotheses**:  
 $H = \{h | h : \mathbb{X} \rightarrow \mathbb{Y}, h(\mathbf{x}) = \text{sign}(\sum_{d=1}^D w_d x_d + b)\}.$
- Learning goal:  $\hat{y} = h(\mathbf{x})$ 
  - ▶ Learn  $w_1, \dots, w_D, b$ .
  - ▶ **Parameters**:  $w_1, \dots, w_D, b$ .
  - ▶  $\mathbf{w}$ : **weights**,  $b$ : **bias**

# Perceptron predict

- Input:  $\mathbf{x} \in \mathbb{R}^D$ ,  $\mathbf{w} \in \mathbb{R}^D$ ,  $b \in \mathbb{R}$ .

$$a = \sum_{d=1}^D w_d x_d + b = \mathbf{w}^T \mathbf{x} + b$$
$$\hat{y} = \text{sign}(a)$$

- Output:  $\hat{y}$ .
- $\sum_{d=1}^D w_d x_d + b = \mathbf{w}^T \mathbf{x} + b = 0$  : hyperplane in  $D$  dimensions with parameters  $(\mathbf{w}, b)$ .
- $\mathbf{w}$ : weights,  $b$ : bias
- $a$ : activation

## Hyperplanes through the origin

Consider  $\mathbf{x}$  that satisfies  $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = 0$ . These  $\mathbf{x}$  define a hyperplane in  $D$  dimensions.

We can always write this as a hyperplane passing through the origin in  $D + 1$  dimensions.

$$\begin{aligned}\tilde{\mathbf{x}} &\equiv \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_D \end{pmatrix} \quad \tilde{\mathbf{w}} \equiv \begin{pmatrix} b \\ w_1 \\ \vdots \\ w_D \end{pmatrix} \\ \tilde{g}(\tilde{\mathbf{x}}) &= \tilde{\mathbf{w}}^T \tilde{\mathbf{x}} \\ &= \sum_{d=1}^D w_d x_d + b \\ &= g(\mathbf{x})\end{aligned}$$



# Perceptron learning

**If we have only one training example  $(\mathbf{x}_n, y_n)$ .**

Assume  $b = 0$ .

How can we change  $\mathbf{w}$  such that

$$y_n = \text{sign}(\mathbf{w}^T \mathbf{x}_n)$$

## Two cases

- If  $y_n = \text{sign}(\mathbf{w}^T \mathbf{x}_n)$ , do nothing.
- If  $y_n \neq \text{sign}(\mathbf{w}^T \mathbf{x}_n)$ ,

$$\mathbf{w}^{\text{NEW}} \leftarrow \mathbf{w}^{\text{OLD}} + y_n \mathbf{x}_n$$

# Perceptron learning

**If we have only one training example  $(\mathbf{x}_n, y_n)$ .**

Assume  $b = 0$ .

How can we change  $\mathbf{w}$  such that

$$y_n = \text{sign}(\mathbf{w}^T \mathbf{x}_n)$$

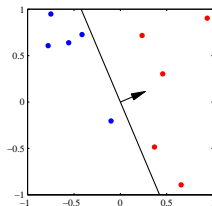
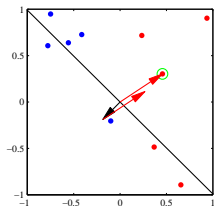
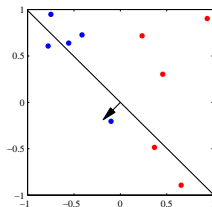
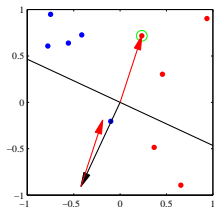
**Another way of saying the same thing**

- $a = \mathbf{w}^T \mathbf{x}_n$
- If  $y_n a > 0$ , do nothing.
- If  $y_n a \leq 0$ ,

$$\mathbf{w}^{\text{NEW}} \leftarrow \mathbf{w}^{\text{OLD}} + y_n \mathbf{x}_n$$

# Example of perceptron update

Red is  $+1$ , Blue is  $-1$



# Why would it work?

If  $y_n a \leq 0$ , then

$$y_n(\mathbf{w}^T \mathbf{x}_n) \leq 0$$

# Why would it work?

If  $y_n a \leq 0$ , then

$$y_n(\mathbf{w}^T \mathbf{x}_n) \leq 0$$

What would happen if we change to new  $\mathbf{w}^{\text{NEW}} = \mathbf{w} + y_n \mathbf{x}_n$ ?

$$y_n[(\mathbf{w} + y_n \mathbf{x}_n)^T \mathbf{x}_n] = y_n \mathbf{w}^T \mathbf{x}_n + y_n^2 \mathbf{x}_n^T \mathbf{x}_n$$

# Why would it work?

If  $y_n a \leq 0$ , then

$$y_n(\mathbf{w}^T \mathbf{x}_n) \leq 0$$

What would happen if we change to new  $\mathbf{w}^{\text{NEW}} = \mathbf{w} + y_n \mathbf{x}_n$ ?

$$y_n[(\mathbf{w} + y_n \mathbf{x}_n)^T \mathbf{x}_n] = y_n \mathbf{w}^T \mathbf{x}_n + y_n^2 \mathbf{x}_n^T \mathbf{x}_n$$

We are adding a positive number, so it is possible that

$$y_n(\mathbf{w}^{\text{NEW}T} \mathbf{x}_n) > 0$$

i.e., we are more likely to classify correctly

# Perceptron learning

## Iteratively solving one case at a time

- REPEAT
- Pick a data point  $\mathbf{x}_n$
- Compute  $a = \mathbf{w}^T \mathbf{x}_n$  using the *current*  $\mathbf{w}$
- If  $ay_n > 0$ , do nothing. Else,

$$\mathbf{w} \leftarrow \mathbf{w} + y_n \mathbf{x}_n$$

- UNTIL converged.

# Perceptron training/learning

*data* = **N samples/instances**:  $= \{(x_1, y_1), \dots, (x_N, y_N)\}$

---

**Algorithm 1** PerceptronTrain (*data*, *maxIter*)

---

```
1:  $w_d \leftarrow 0, d \in \{0, D\}$ 
2: for  $iter = 1 \dots MaxIter$  do
3:   for  $(x, y) \in data$  do
4:      $a \leftarrow w^T x$ 
5:     if  $ay \leq 0$  then
6:        $w \leftarrow w + yx$ 
7:     end if
8:   end for
9: end for
10: return  $w$ 
```

---



# Design decisions

- *MaxIter*: Hyperparameter

# Design decisions

- *MaxIter*: Hyperparameter
- How to loop over the data?
  - ▶ Constant.
  - ▶ Permuting once
  - ▶ Permuting in each iteration

# Properties of perceptron learning

- This is an **online** algorithm – looks at one instance at a time.
- Does the algorithm terminate (**convergence**)?

# Properties of perceptron learning

- This is an **online** algorithm – looks at one instance at a time.
- Does the algorithm terminate (**convergence**)?
  - ▶ If training data is **not linearly separable**, the algorithm does not converge.
  - ▶ If the training data is linearly separable, the algorithm stops in a finite number of steps (**converges**).

# Properties of perceptron learning

- This is an **online** algorithm – looks at one instance at a time.
- Does the algorithm terminate (**convergence**)?
  - ▶ If training data is **not linearly separable**, the algorithm does not converge.
  - ▶ If the training data is linearly separable, the algorithm stops in a finite number of steps (**converges**).
- How long to convergence ?
  - ▶ Depends on the difficulty of the problem.

# Perceptron

- Extensions
  - ▶ Voting
  - ▶ Averaging
- Limitations
  - ▶ Linear separability
- Interpreting the importance of features
  - ▶ The values of weight  $w_d$  tells us the importance of feature  $x_d$ .

# Outline

- 1 Review of previous lecture
- 2 Perceptron
- 3 What we have learned

# Summary

- You should now be able to understand the differences between decision trees, perceptrons and nearest neighbors.
- Given data, use training, development and test splits (or cross-validation).
- Use training and development to tune hyperparameters that trades off overfitting and underfitting.
- Use test to get an estimate of generalization or accuracy on unseen data.