



SOI1010

Machine Learning II

Lecture 4: Linear Regression Cont'd & Classification

Department of Data Science
Hanyang University

Recap: Machine Learning

- Learn from data to achieve generalization:
 - To design algorithms that can be trained on data to make predictions for new examples/data
- Supervised Learning

Assume there is a mapping function:



Train a model f that is as close to the mapping function as possible

- To make correct predictions for new examples



Recap: Math formulation

- Given training data $\{(\mathbf{x}^{(i)}, y^{(i)}) : 1 \leq i \leq n\}$ i.i.d. from distribution \mathcal{D}
- Find $y = f(\mathbf{x}) \in \mathcal{H}$ that minimizes $\hat{\mathcal{L}}(f) = \frac{1}{n} \sum_{i=1}^n l(f, \mathbf{x}_i, y_i)$
- s.t. the expected loss is small

$$\mathcal{L}(f) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}}[l(f, \mathbf{x}, y)]$$

Recap: Machine Learning 1-2-3

- Collect data and extract features
- Build model: choose hypothesis class \mathcal{H} and loss function l
- Optimization: minimize the empirical loss

Experience

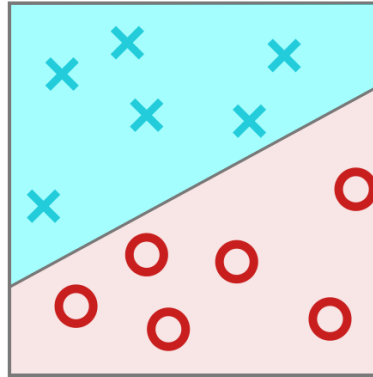
The diagram consists of two yellow-outlined boxes. The 'Experience' box is at the top right, with an arrow pointing to the first bullet point 'Collect data and extract features'. The 'Prior knowledge' box is at the bottom right, with three arrows pointing to the second, third, and fourth bullet points: 'Build model: choose hypothesis class \mathcal{H} and loss function l ', 'Optimization: minimize the empirical loss', and 'Optimization: minimize the empirical loss'.

Prior knowledge

Recap: Classification vs Regression

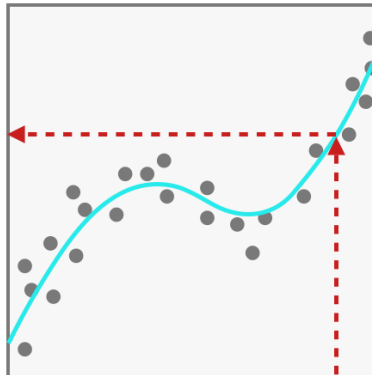
- Classification

- The task of predicting a category



- Regression

- The task of predicting a continuous quantity




Recap: Linear regression math formulation

- Given training data $\{(\mathbf{x}^{(i)}, y^{(i)}) : 1 \leq i \leq n\}$ i.i.d. from distribution \mathcal{D}
- Find $f_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ that minimizes $\hat{\mathcal{L}}(f_{\mathbf{w}}) = \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}^{(i)} - y^{(i)})^2$

Linear model \mathcal{H}



l_2 loss



Why l_2 loss?

- Why not choose another loss
 - l_1 loss, hinge loss, exponential loss, etc...
- Empirical: easy to optimize
 - For linear case: $\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$
- Theoretical: a way to encode prior knowledge
- Questions:
 - What kind of prior knowledge?
 - Principal way to derive loss function?

Maximum likelihood Estimation (MLE)

Maximum likelihood Estimation (MLE)

- Given training data $\{(\mathbf{x}^{(i)}, y^{(i)}) : 1 \leq i \leq n\}$ i.i.d. from distribution \mathcal{D}
- Let $\{P_{\theta}(\mathbf{x}, y) : \theta \in \Theta\}$ be a family of distributions parameterized by θ
- Would like to pick θ so that $P_{\theta}(\mathbf{x}, y)$ fits the data well

Maximum likelihood Estimation (MLE)

- Given training data $\{(\mathbf{x}^{(i)}, y^{(i)}) : 1 \leq i \leq n\}$ i.i.d. from distribution \mathcal{D}
- Let $\{P_{\theta}(\mathbf{x}, y) : \theta \in \Theta\}$ be a family of distributions parameterized by θ
- “fitness” of θ to one data point $(\mathbf{x}^{(i)}, y^{(i)})$

$$\text{likelihood}(\theta; \mathbf{x}^{(i)}, y^{(i)}) := P_{\theta}(\mathbf{x}^{(i)}, y^{(i)})$$

Maximum likelihood Estimation (MLE)

- Given training data $\{(\mathbf{x}^{(i)}, y^{(i)}) : 1 \leq i \leq n\}$ i.i.d. from distribution \mathcal{D}
- Let $\{P_{\theta}(\mathbf{x}, y) : \theta \in \Theta\}$ be a family of distributions parameterized by θ
- “fitness” of θ to **i.i.d.** data points $\{(\mathbf{x}^{(i)}, y^{(i)})\}$

$$\text{likelihood}(\theta; \{\mathbf{x}^{(i)}, y^{(i)}\}) := P_{\theta}(\{\mathbf{x}^{(i)}, y^{(i)}\}) = \prod_i P_{\theta}(\mathbf{x}^{(i)}, y^{(i)})$$

Maximum likelihood Estimation (MLE)

- Given training data $\{(\mathbf{x}^{(i)}, y^{(i)}) : 1 \leq i \leq n\}$ i.i.d. from distribution \mathcal{D}
- Let $\{P_{\boldsymbol{\theta}}(\mathbf{x}, y) : \boldsymbol{\theta} \in \Theta\}$ be a family of distributions parameterized by $\boldsymbol{\theta}$
- **MLE: maximize** “fitness” of $\boldsymbol{\theta}$ to i.i.d. data points $\{(\mathbf{x}^{(i)}, y^{(i)})\}$

$$\begin{aligned}\boldsymbol{\theta}_{\text{ML}} &= \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \text{likelihood}(\boldsymbol{\theta}; \{\mathbf{x}^{(i)}, y^{(i)}\}) \\ &= \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \prod_i P_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}, y^{(i)})\end{aligned}$$

Maximum likelihood Estimation (MLE)

- Given training data $\{(\mathbf{x}^{(i)}, y^{(i)}) : 1 \leq i \leq n\}$ i.i.d. from distribution \mathcal{D}
- Let $\{P_{\boldsymbol{\theta}}(\mathbf{x}, y) : \boldsymbol{\theta} \in \Theta\}$ be a family of distributions parameterized by $\boldsymbol{\theta}$
- **MLE: maximize** “fitness” of $\boldsymbol{\theta}$ to i.i.d. data points $\{(\mathbf{x}^{(i)}, y^{(i)})\}$

$$\boldsymbol{\theta}_{\text{ML}} = \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \log \left[\prod_i P_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}, y^{(i)}) \right]$$

$$\boldsymbol{\theta}_{\text{ML}} = \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \sum_i \log [P_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}, y^{(i)})]$$

Maximum likelihood Estimation (MLE)

- Given training data $\{(\mathbf{x}^{(i)}, y^{(i)}) : 1 \leq i \leq n\}$ i.i.d. from distribution \mathcal{D}
- Let $\{P_{\boldsymbol{\theta}}(\mathbf{x}, y) : \boldsymbol{\theta} \in \Theta\}$ be a family of distributions parameterized by $\boldsymbol{\theta}$
- MLE: negative log-likelihood loss

$$\boldsymbol{\theta}_{\text{ML}} = \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \sum_i \log[P_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}, y^{(i)})]$$

$$l(P_{\boldsymbol{\theta}}, \mathbf{x}^{(i)}, y^{(i)}) = -\log(P_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}, y^{(i)}))$$

$$\hat{\mathcal{L}}(P_{\boldsymbol{\theta}}) = -\sum_i \log(P_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}, y^{(i)}))$$

MLE: conditional log-likelihood

- Given training data $\{(\mathbf{x}^{(i)}, y^{(i)}) : 1 \leq i \leq n\}$ i.i.d. from distribution \mathcal{D}
- Let $\{P_{\boldsymbol{\theta}}(y|\mathbf{x}) : \boldsymbol{\theta} \in \Theta\}$ be a family of distributions parameterized by $\boldsymbol{\theta}$
- MLE: negative **conditional** log-likelihood loss

Only care about predicting y from \mathbf{x} ; do not care about $P(\mathbf{x}, y)$

$$\boldsymbol{\theta}_{\text{ML}} = \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \sum_i \log[P_{\boldsymbol{\theta}}(y^{(i)}|\mathbf{x}^{(i)})]$$

$$l(P_{\boldsymbol{\theta}}, \mathbf{x}^{(i)}, y^{(i)}) = -\log(P_{\boldsymbol{\theta}}(y^{(i)}|\mathbf{x}^{(i)}))$$

$$\hat{\mathcal{L}}(P_{\boldsymbol{\theta}}) = -\sum_i \log(P_{\boldsymbol{\theta}}(y^{(i)}|\mathbf{x}^{(i)}))$$

MLE: conditional log-likelihood

- Given training data $\{(\mathbf{x}^{(i)}, y^{(i)}) : 1 \leq i \leq n\}$ i.i.d. from distribution \mathcal{D}
- Let $\{P_{\boldsymbol{\theta}}(y|\mathbf{x}) : \boldsymbol{\theta} \in \Theta\}$ be a family of distributions parameterized by $\boldsymbol{\theta}$
- MLE: negative **conditional** log-likelihood loss

$P(y|\mathbf{x})$: discriminative;
 $P(\mathbf{x}, y)$: generative

$$\boldsymbol{\theta}_{\text{ML}} = \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \sum_i \log[P_{\boldsymbol{\theta}}(y^{(i)}|\mathbf{x}^{(i)})]$$

$$l(P_{\boldsymbol{\theta}}, \mathbf{x}^{(i)}, y^{(i)}) = -\log(P_{\boldsymbol{\theta}}(y^{(i)}|\mathbf{x}^{(i)}))$$

$$\hat{\mathcal{L}}(P_{\boldsymbol{\theta}}) = -\sum_i \log(P_{\boldsymbol{\theta}}(y^{(i)}|\mathbf{x}^{(i)}))$$

MLE: conditional log-likelihood

- Given training data $\{(\mathbf{x}^{(i)}, y^{(i)}) : 1 \leq i \leq n\}$ i.i.d. from distribution \mathcal{D}
- Let $\{P_{\boldsymbol{\theta}}(y|\mathbf{x}) : \boldsymbol{\theta} \in \Theta\}$ be a family of distributions parameterized by $\boldsymbol{\theta}$
- MLE: negative **conditional** log-likelihood loss

$P(y|\mathbf{x})$: discriminative;
 $P(\mathbf{x}, y)$: generative

$$\boldsymbol{\theta}_{\text{ML}} = \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \sum_i \log[P_{\boldsymbol{\theta}}(y^{(i)}|\mathbf{x}^{(i)})]$$

$$l(P_{\boldsymbol{\theta}}, \mathbf{x}^{(i)}, y^{(i)}) = -\log(P_{\boldsymbol{\theta}}(y^{(i)}|\mathbf{x}^{(i)}))$$

$$\hat{\mathcal{L}}(P_{\boldsymbol{\theta}}) = -\sum_i \log(P_{\boldsymbol{\theta}}(y^{(i)}|\mathbf{x}^{(i)}))$$

MLE: conditional log-likelihood

- Given training data $\{(\mathbf{x}^{(i)}, y^{(i)}) : 1 \leq i \leq n\}$ i.i.d. from distribution \mathcal{D}
- Let $\{P_{\boldsymbol{\theta}}(y|\mathbf{x}) : \boldsymbol{\theta} \in \Theta\}$ be a family of distributions parameterized by $\boldsymbol{\theta}$
- Define $P_{\boldsymbol{\theta}}(y|\mathbf{x}) = \mathcal{N}(y; f_{\boldsymbol{\theta}}(\mathbf{x}), \sigma^2)$

$$\log(P_{\boldsymbol{\theta}}(y^{(i)}|\mathbf{x}^{(i)})) = \frac{-1}{2\sigma^2}(f_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y_i)^2 - \log(\sigma) - \frac{1}{2}\log(2\pi)$$

$$\boldsymbol{\theta}_{\text{ML}} = \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \sum_i \log[P_{\boldsymbol{\theta}}(y^{(i)}|\mathbf{x}^{(i)})]$$

$$\hat{\mathcal{L}}(P_{\boldsymbol{\theta}}) = - \sum_i \log(P_{\boldsymbol{\theta}}(y^{(i)}|\mathbf{x}^{(i)}))$$

$$\boldsymbol{\theta}_{\text{ML}} = \operatorname{argmin}_{\boldsymbol{\theta} \in \Theta} \hat{\mathcal{L}}(P_{\boldsymbol{\theta}}) = \operatorname{argmin}_{\boldsymbol{\theta} \in \Theta} \frac{1}{n} \sum_{i=1}^n (f_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)})^2$$

MLE: conditional log-likelihood

- Given training data $\{(\mathbf{x}^{(i)}, y^{(i)}) : 1 \leq i \leq n\}$ i.i.d. from distribution \mathcal{D}
- Find $f_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ that minimizes $\hat{\mathcal{L}}(f_{\mathbf{w}}) = \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}^{(i)} - y^{(i)})^2$
- Define $P_{\boldsymbol{\theta}}(y|\mathbf{x}) = \mathcal{N}(y; f_{\boldsymbol{\theta}}(\mathbf{x}), \sigma^2)$

$$\log(P_{\boldsymbol{\theta}}(y^{(i)}|\mathbf{x}^{(i)})) = \frac{-1}{2\sigma^2} (f_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y_i)^2 - \log(\sigma) - \frac{1}{2} \log(2\pi)$$

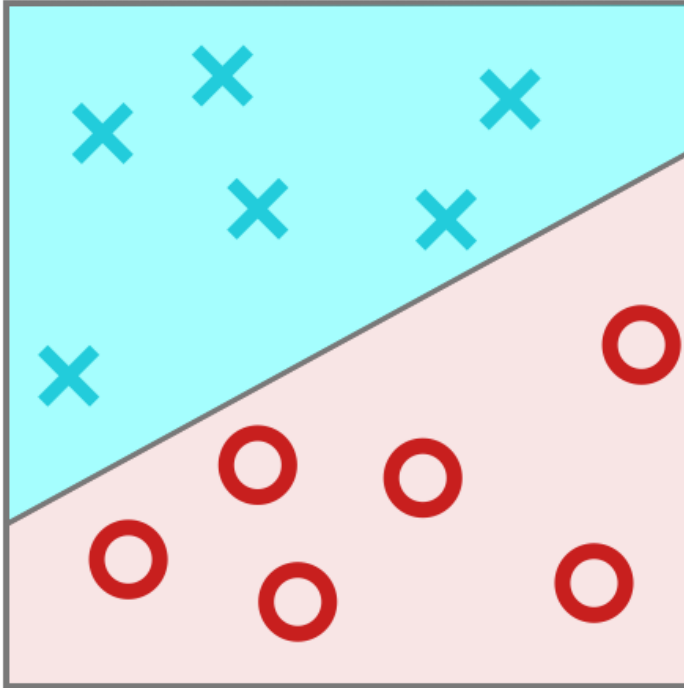
$$\boldsymbol{\theta}_{\text{ML}} = \operatorname{argmax}_{\boldsymbol{\theta} \in \Theta} \sum_i \log[P_{\boldsymbol{\theta}}(y^{(i)}|\mathbf{x}^{(i)})]$$

l_2 loss: Normal + MLE

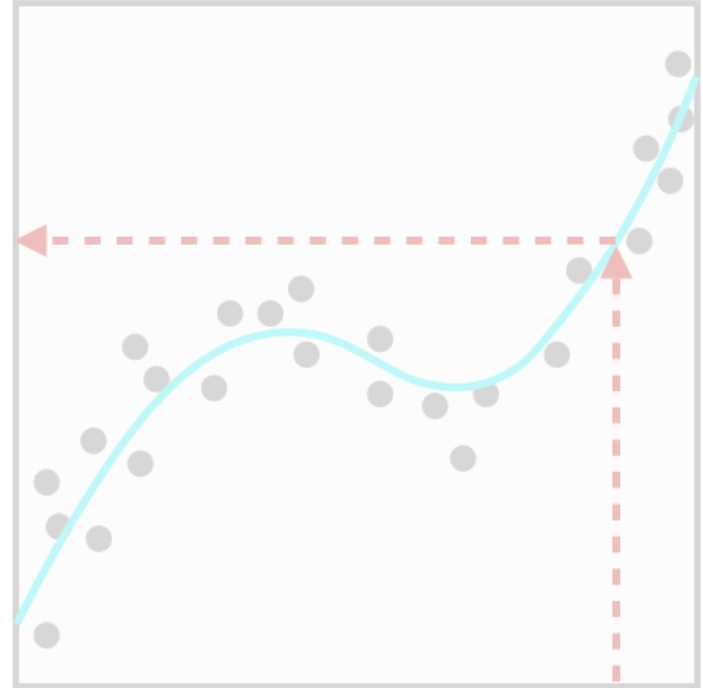
$$\hat{\mathcal{L}}(P_{\boldsymbol{\theta}}) = - \sum_i \log(P_{\boldsymbol{\theta}}(y^{(i)}|\mathbf{x}^{(i)}))$$

$$\boldsymbol{\theta}_{\text{ML}} = \operatorname{argmin}_{\boldsymbol{\theta} \in \Theta} \hat{\mathcal{L}}(P_{\boldsymbol{\theta}}) = \operatorname{argmin}_{\boldsymbol{\theta} \in \Theta} \frac{1}{n} \sum_{i=1}^n (f_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)})^2$$

Classification vs Regression



Classification



Regression

Classification

Why classification?

- Classification: a kind of summary
- Easy to interpret
- Easy for making decisions
- Example: image classification

Computer Vision

- Computer Vision

- Subfield of machine learning
- Goal: build the computer systems that can understand images and/or videos the way humans do

- What do you see in the image?

- What is he doing?



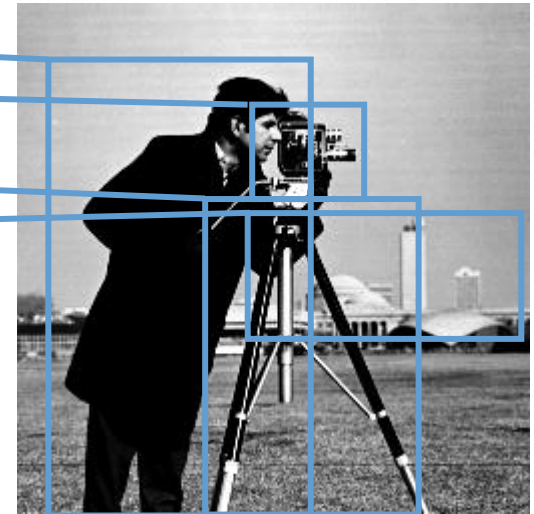
Computer Vision

- Computer Vision

- Subfield of machine learning
- Goal: build the computer systems that can understand images and/or videos the way humans do

- What do you see in the image?

- Person
- Camera
- Tripod
- buildings



- What is he doing?

Computer Vision

- Computer Vision

- Subfield of machine learning
- Goal: build the computer systems that can understand images and/or videos the way humans do

- What do you see in the image?

- What is he doing?

- Taking a picture of something
- Or maybe taking a pose...?



Computer Vision

- Computer Vision

- Subfield of machine learning
- Goal: build the computer systems that can understand images and/or videos the way humans do

- What do computers see in the image?



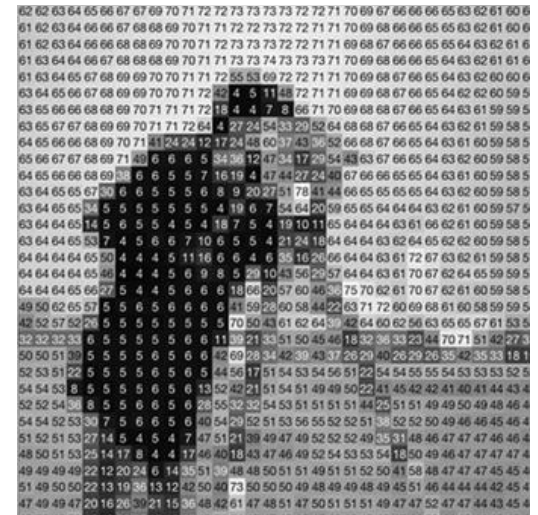
Computer Vision

- Computer Vision

- Subfield of machine learning
- Goal: build the computer systems that can understand images and/or videos the way humans do

- What do computers see in the image?

- Just numbers..!



Computer Vision

- Computer Vision

- Subfield of machine learning
- Goal: build the computer systems that can understand images and/or videos the way humans do

- How can we bridge the gap?

02 02 03 04 05 06 07 07 09 70 71 72 72 73 73 73 73 72 72 71 70 09 07 06 06 06 05 03 02 01 00 0
01 02 03 04 06 06 07 08 08 09 70 71 71 72 72 73 72 72 71 71 70 09 08 06 06 05 05 03 02 01 00 0
01 02 03 04 06 06 08 08 09 70 71 72 73 73 72 72 71 71 09 08 07 06 06 05 05 04 03 02 01 0
01 03 04 04 06 07 08 08 09 70 71 71 73 74 73 73 73 71 70 09 08 06 06 05 04 03 02 01 01 0
01 03 04 05 07 08 09 09 70 71 71 72 55 09 72 72 71 71 70 09 08 07 06 05 04 03 02 00 00 0
03 04 05 06 07 08 09 09 70 71 72 02 4 5 11 08 72 71 71 09 09 08 07 06 05 04 02 02 00 59 5
03 05 06 06 08 09 70 71 71 72 18 4 4 7 8 06 71 70 09 08 07 06 05 04 03 01 59 59 5
03 05 07 07 08 09 09 70 71 71 72 04 4 27 24 54 33 29 52 04 08 08 07 06 05 04 03 02 01 59 58 5
04 05 06 06 08 09 70 71 71 72 04 12 17 24 43 00 37 43 06 52 06 08 07 06 05 04 03 01 00 59 58 5
05 06 07 07 08 09 71 00 6 6 6 5 84 08 12 47 34 17 29 54 43 03 07 06 05 04 03 02 00 59 58 5
04 05 06 06 08 09 00 6 6 5 5 7 16 19 4 47 44 27 24 40 07 06 06 05 05 04 03 01 00 59 58 5
03 04 05 05 07 00 6 6 5 5 5 6 8 9 20 27 51 78 11 44 06 05 05 05 05 04 03 02 00 59 58 5
03 04 05 05 04 5 5 5 5 5 5 5 4 19 6 7 54 64 20 59 05 05 04 04 04 03 02 01 00 59 57 5
03 04 04 05 14 5 6 5 5 4 5 4 18 7 5 4 19 10 11 05 04 04 04 03 01 06 02 01 00 59 58 5
03 04 04 05 03 7 4 5 6 6 7 10 6 5 5 4 21 24 18 04 04 04 03 02 04 05 02 02 00 59 58 5
04 04 04 04 05 50 4 4 4 5 11 16 6 6 4 6 35 16 26 06 04 04 03 01 72 07 03 02 01 59 58 5
04 04 04 04 05 46 4 4 4 5 6 9 8 5 29 10 43 56 29 57 04 04 03 01 70 07 02 04 05 59 59 5
04 04 04 05 06 27 5 4 4 5 6 6 6 18 06 20 57 60 46 03 75 70 02 01 70 07 02 01 00 59 58 5
49 50 02 05 57 5 5 5 6 6 6 6 6 41 59 26 00 58 44 22 03 71 72 00 09 08 01 00 58 59 59 5
42 52 57 52 26 5 5 5 5 5 5 5 5 70 50 43 61 02 64 39 42 64 00 02 56 03 05 05 07 01 03 6
32 32 32 33 6 5 5 5 5 5 6 6 11 39 21 33 51 50 45 46 18 32 36 33 23 44 70 71 51 42 27 5
50 50 51 09 5 5 5 5 5 5 6 6 42 09 28 34 42 39 43 37 26 29 40 26 29 26 35 42 35 33 18 1
52 53 51 22 5 5 5 5 5 5 6 6 54 56 17 51 54 53 54 56 51 22 54 54 55 55 54 53 53 53 52 5
54 54 53 8 5 5 5 5 5 6 6 13 52 42 21 51 54 51 49 49 50 22 41 45 42 42 41 40 41 44 43 4
52 52 54 50 8 5 5 6 6 5 6 28 55 32 32 54 53 51 51 51 51 44 25 51 51 49 49 50 49 48 46 4
54 54 52 53 30 7 5 6 6 5 6 40 54 29 52 51 53 56 55 52 52 51 38 52 52 50 49 46 46 45 46 4
51 52 51 53 27 14 5 4 5 4 7 47 51 21 39 49 47 49 52 52 52 49 53 31 48 46 47 47 47 46 46 4
48 50 51 53 25 14 17 8 4 4 17 46 40 18 43 47 46 49 52 54 53 53 54 18 50 49 46 47 47 47 47 4
49 49 49 49 22 12 20 24 6 14 35 51 39 48 48 50 51 51 49 51 51 52 50 41 58 48 47 47 47 45 45 4
51 49 50 50 22 13 19 36 13 12 42 50 40 73 50 50 50 49 48 49 48 49 45 51 46 44 44 44 42 45 4
47 49 49 47 20 16 26 30 21 15 36 48 42 61 47 48 51 47 50 51 51 51 49 47 47 52 47 47 44 43 45 4



Computer Vision

- Computer Vision

- Subfield of machine learning
- Goal: build the computer systems that can understand images and/or videos the way humans do

- How can we bridge the gap?



Computer Vision

- To understand/predict behaviors of objects/targets
(Action recognition/prediction)
- We need to understand their interactions with objects and surroundings
(Scene understanding)
- We need to understand the locations of objects
(Localization)
- We need to understand what the objects are
(Classification)

Image Classification

- One of the most core tasks in computer vision
- Binary classification
 - Whether the image has an object of interest
- e.g.,



This image by Nikita is
licensed under [CC-BY 2.0](#)



cat

Image Classification

There are several challenges

- Viewpoint variation



Image Classification

There are several challenges

- Background clutter



[This image](#) is [CC0 1.0](#) public domain



[This image](#) is [CC0 1.0](#) public domain

Image Classification

There are several challenges

- Illumination



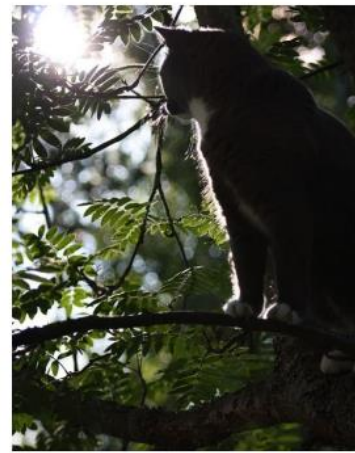
[This image](#) is [CC0 1.0](#) public domain



[This image](#) is [CC0 1.0](#) public domain



[This image](#) is [CC0 1.0](#) public domain



[This image](#) is [CC0 1.0](#) public domain

Image Classification

There are several challenges

- Occlusion



[This image](#) is [CC0 1.0](#) public domain



[This image](#) is [CC0 1.0](#) public domain



[This image](#) by [jonsson](#) is licensed under [CC-BY 2.0](#)

Image Classification

There are several challenges

- Deformation



[This image](#) by [Umberto Salvagnin](#) is licensed under [CC-BY 2.0](#)



[This image](#) by [Umberto Salvagnin](#) is licensed under [CC-BY 2.0](#)



[This image](#) by [sare bear](#) is licensed under [CC-BY 2.0](#)



[This image](#) by [Tom Thai](#) is licensed under [CC-BY 2.0](#)

Image Classification

There are several challenges

- Intraclass variation



[This image](#) is [CC0 1.0](#) public domain.

Image Classification

There are several challenges

- Scale variation



Image Classification

How to solve such a challenging problem?

Viewpoint variation



Scale Variation



Occlusion



Illumination



Deformation



Intra-class



Background clutter



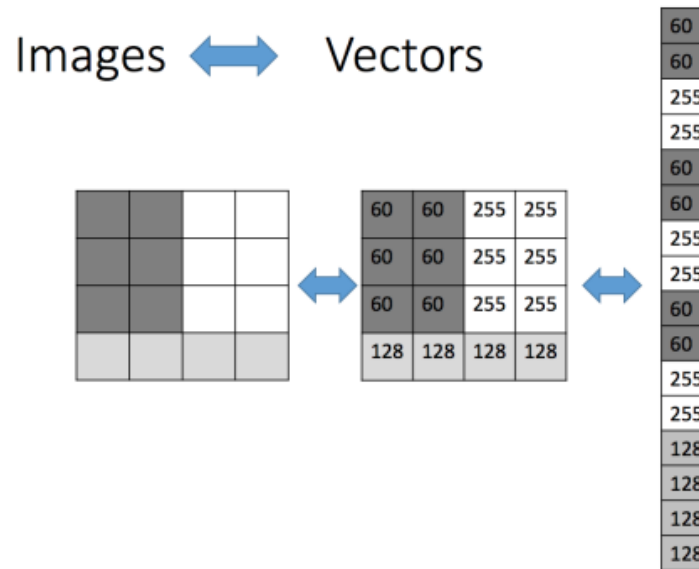
Machine Learning

Representation

- Represent the input as an input vector in \mathbb{R}^d
- Map the data to another space that is easy to manipulate
- Vectors are a great representation since we can do linear algebra!

Machine Learning

Representation



Can do much better if you compute a vector of meaningful features

Later, we will see that Deep Learning is good at obtaining very strong features

Machine Learning

Machine learning algorithms can be classified, depending on the assumptions w.r.t. the form of a function f

- Non-parametric
- Parametric

Machine Learning

Machine learning algorithms can be classified, depending on the assumptions w.r.t. the form of a function f

- Non-parametric

- No assumption on the form of the function
- Free to learn any functional form
- Good when you have a lot of data and no prior knowledge
- The number of parameters grows with the amount of training data
- e.g., k-nearest neighbors, decision tree, random forest

- Parametric

Machine Learning

Machine learning algorithms can be classified, depending on the assumptions w.r.t. the form of a function f

- Non-parametric
- Parametric
 - Assumption with regards to the form of the function
 - Certain hypothesis space
 - Can simplify the learning process, but can also limit what can be learned
 - The model has a fixed number of parameters, regardless of the amount of data
 - e.g., logistic regression, LDA, neural networks

Machine Learning

Machine learning algorithms can be classified, depending on the assumptions w.r.t. the form of a function f

- Non-parametric

- No assumption on the form of the function
- Free to learn any functional form
- Good when you have a lot of data and no prior knowledge
- The number of parameters grows with the amount of training data
- e.g., k-nearest neighbors, decision tree, random forest

- Parametric

Nearest Neighbors

Suppose we're given a novel input vector \mathbf{x} we'd like to classify

Idea: Predict the label as that of the most similar training image

- “most similar”: the nearest input vector to \mathbf{x} in the training set
- Can formalize “nearest” in terms of Euclidean distance

$$\text{dist}(\mathbf{x}^{(a)}, \mathbf{x}^{(b)}) = \|\mathbf{x}^{(a)} - \mathbf{x}^{(b)}\|_2 = \sqrt{\sum_{j=1}^d (x_j^{(a)} - x_j^{(b)})^2}$$

Note: we don't need to compute the square root. Why?

Nearest Neighbors

Algorithm:

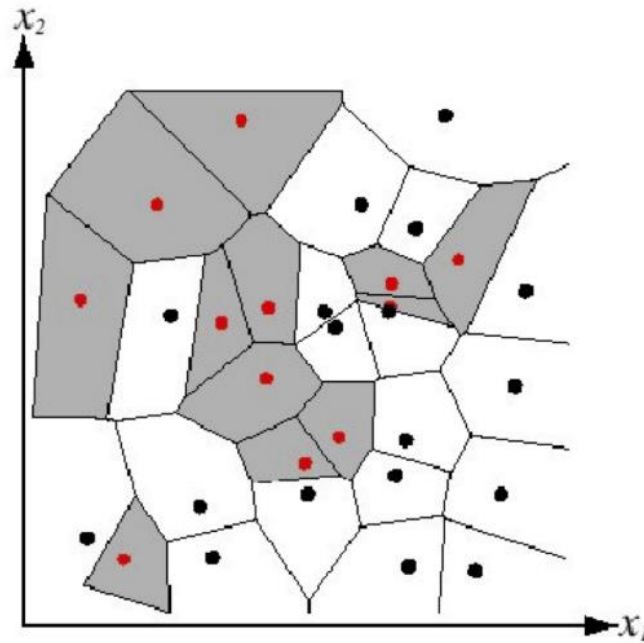
1. Find a training example $(\mathbf{x}^{(*)}, y^{(*)})$ closest to \mathbf{x} :

$$\mathbf{x}^{(*)} = \arg \min_{\mathbf{x}^{(i)} \in \text{train. set}} \text{dist}(\mathbf{x}^{(i)}, \mathbf{x})$$

2. Output $y = y^{*}$

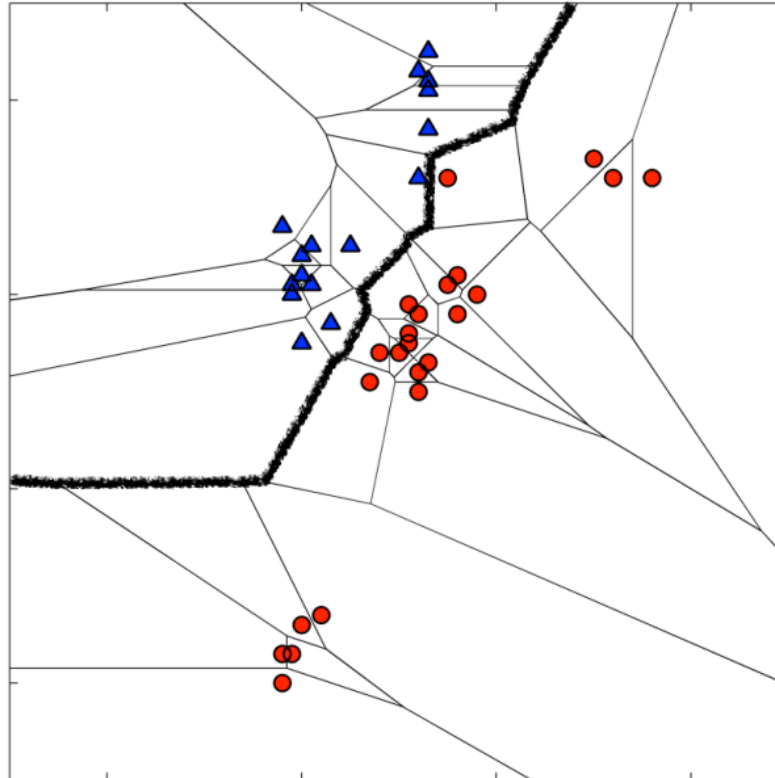
Nearest Neighbors: Decision Boundaries

We can visualize the behavior using a **Voronoi diagram**.



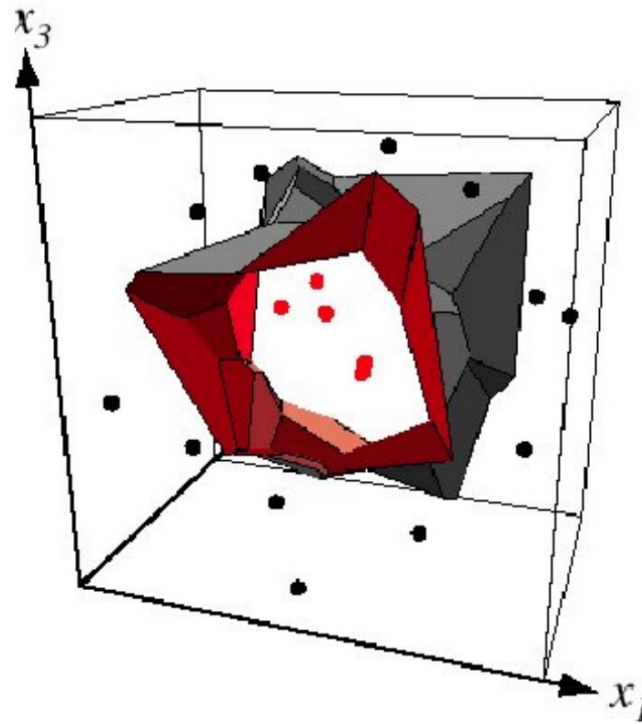
Nearest Neighbors: Decision Boundaries

Decision boundary: the boundary between regions of input space assigned to different categories (classes)



Nearest Neighbors: Decision Boundaries

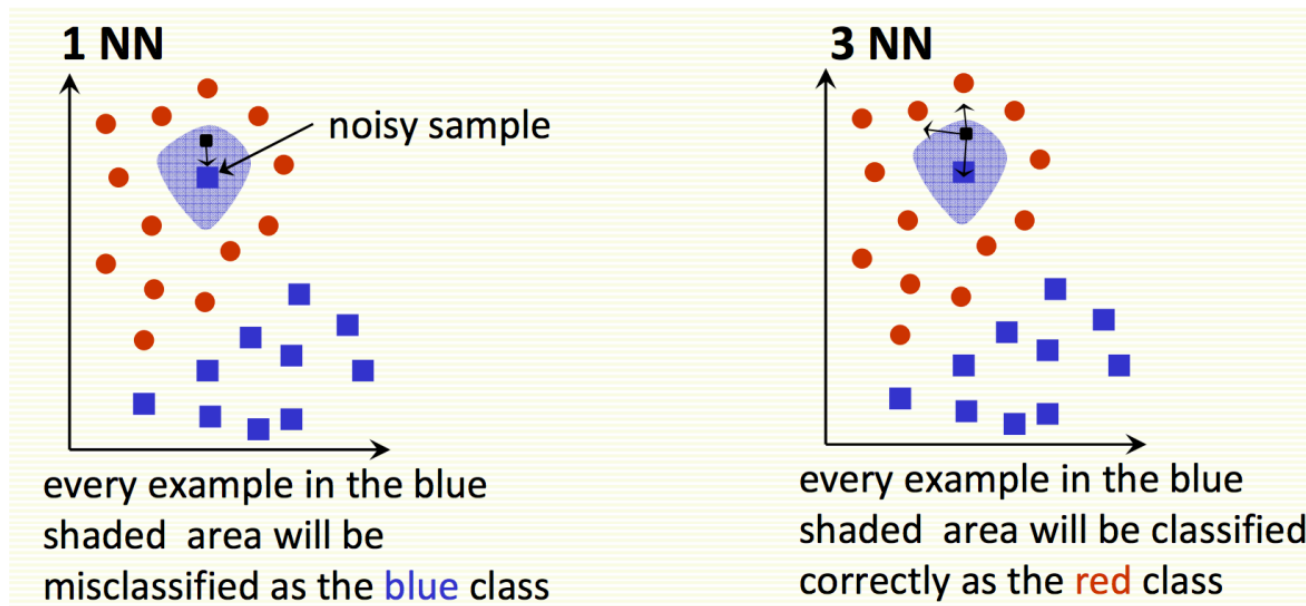
Decision boundary: the boundary between regions of input space assigned to different categories (classes)



k-Nearest Neighbors

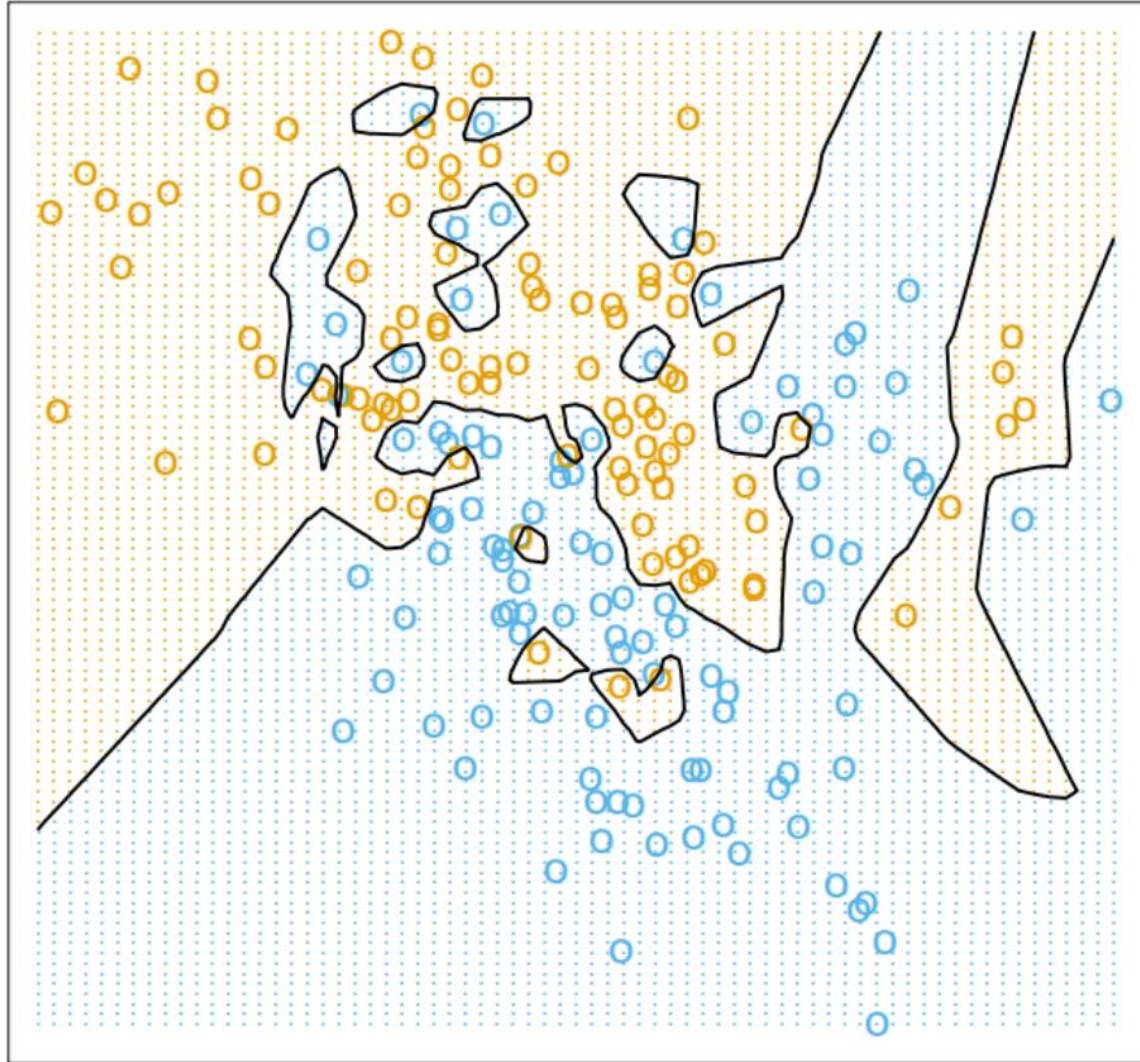
Nearest Neighbors are sensitive to **noise or mis-labeled data**

Take a majority vote from k closest points



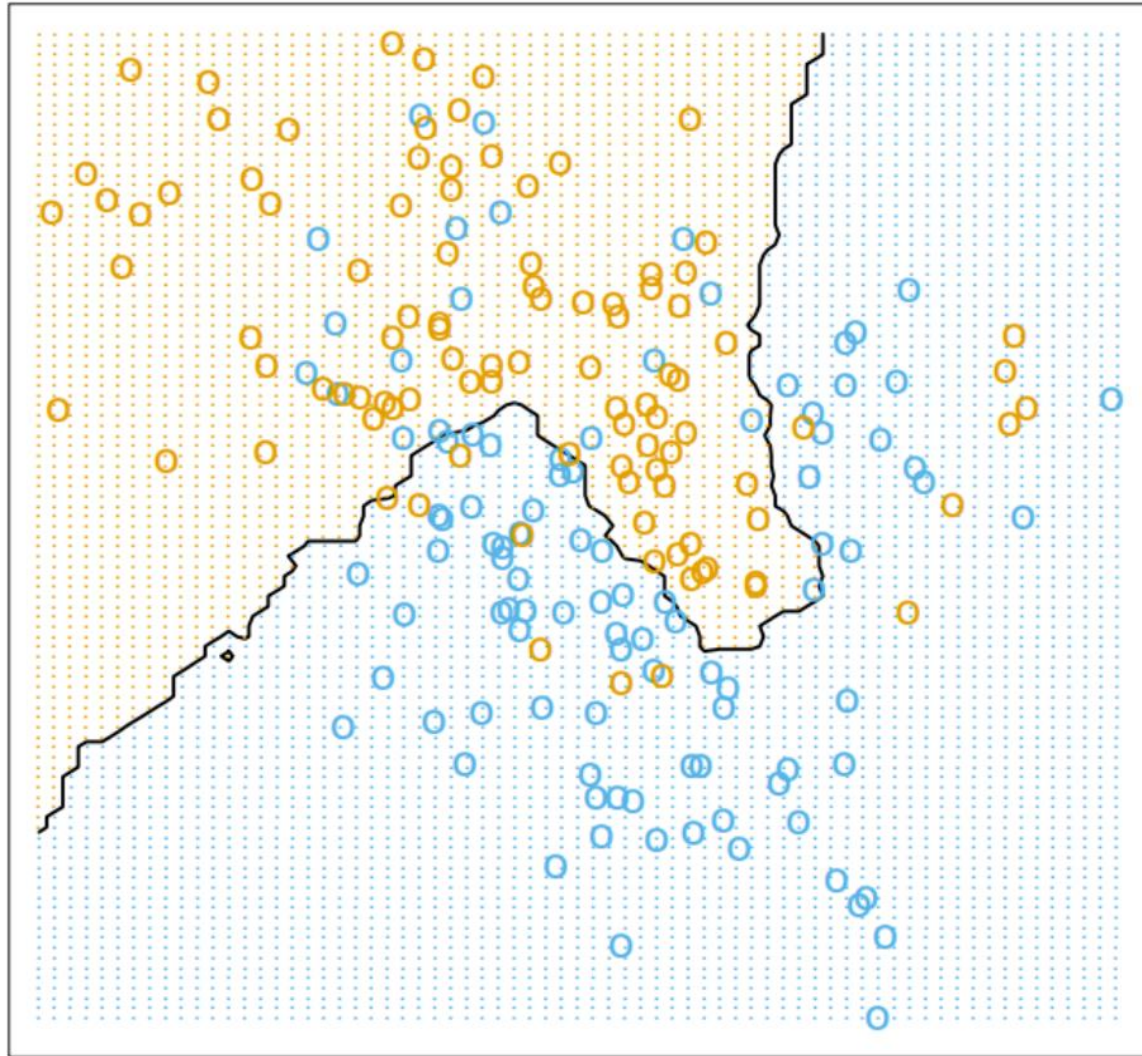
k-Nearest Neighbors

$k=1$



k-Nearest Neighbors

k=15



k-Nearest Neighbors

What is the best value of k ? Tradeoffs in choosing k ?

- Small k
 - Good at capturing fine-grained patterns
 - May **overfit**, i.e., be sensitive to random noise in the training data
- Large k
 - Makes stable predictions by averaging over examples
 - May **underfit**, i.e., fail to capture important information
- **Rule of thumb:** $k < \sqrt{n}$, where n is the number of training examples

k-Nearest Neighbors : Hyperparameters

Design choices about the algorithms themselves

k is a hyperparameter.

Distance measure is a hyperparameter.

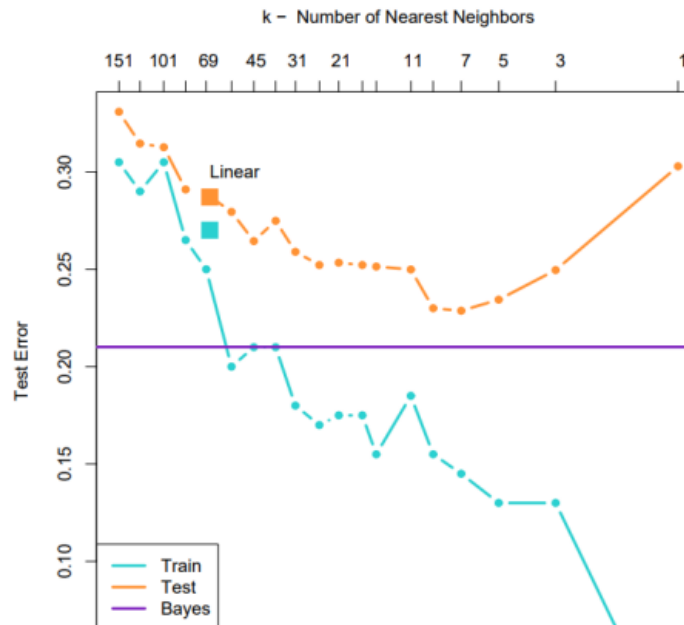
Very problem/dataset-dependent.

Must try them all out and see what works best.

k-Nearest Neighbors : Hyperparameters

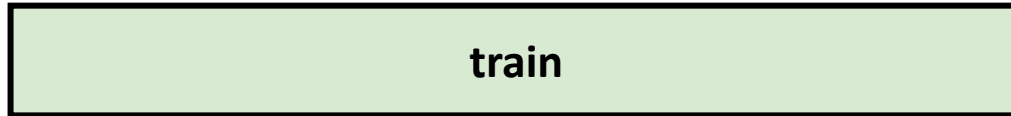
We would like our algorithm to generalize to (perform well on) data it hasn't seen before.

We want to select hyperparameters that give the lowest generalization error (error rate on new examples).



Setting Hyperparameters

Idea #1: Choose hyperparameters that work best on the **training data**



BAD: $k=1$ is the best option for training data

Idea #2: Choose hyperparameters that work best on **test data**

BAD: hyperparameters are overfitted to test data!
No idea how an algorithm will perform on new data



Idea #3: Split **train** data into **train & val**.

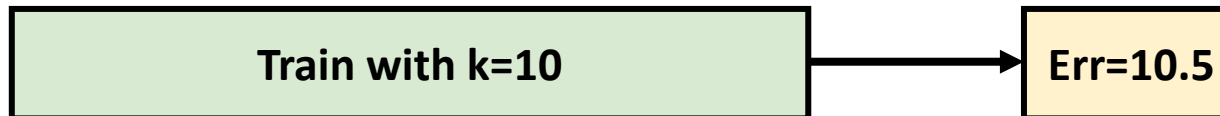
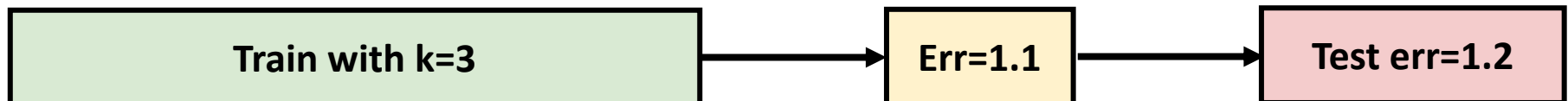
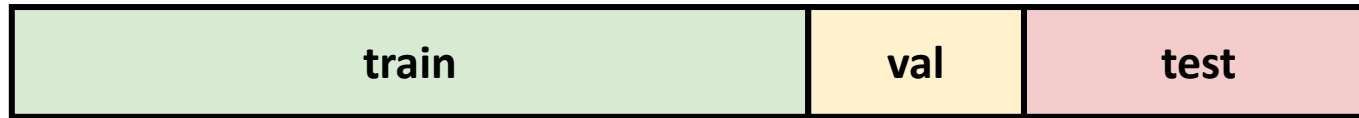
Choose hyperparameters on val and evaluate on test



Setting Hyperparameters

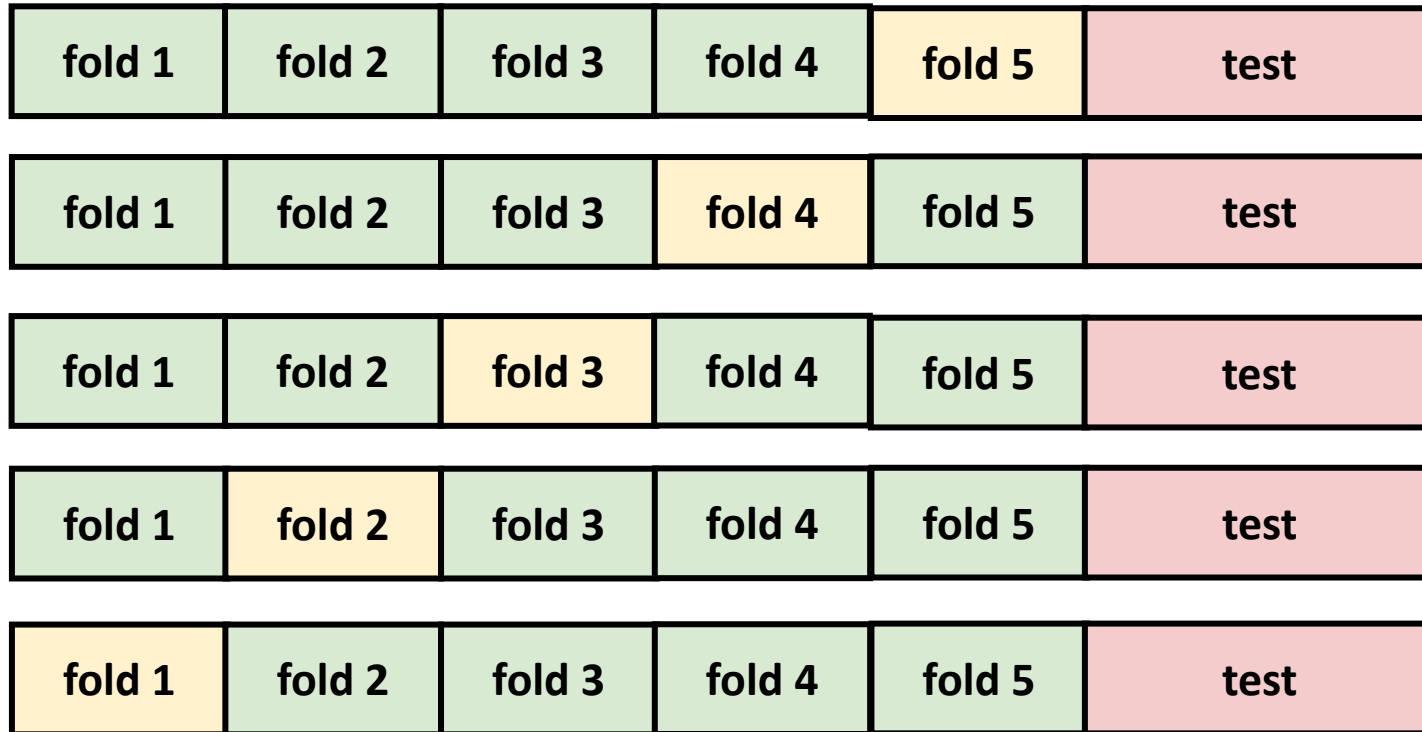
Idea #3: Split **train** data into **train & val**.

Choose hyperparameters on val and evaluate on test



Setting Hyperparameters

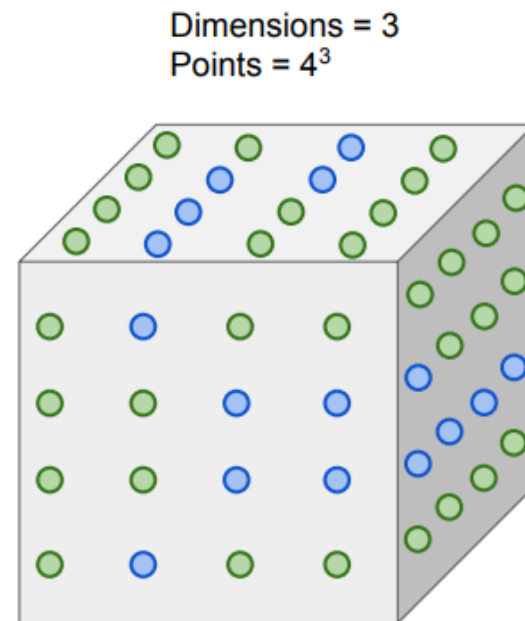
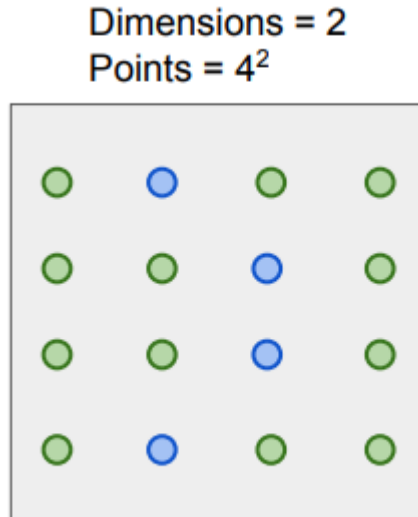
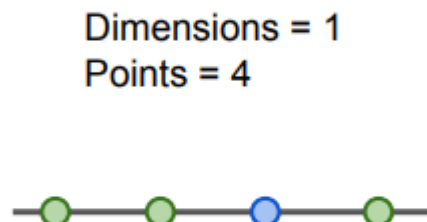
Idea #4: Cross-Validation: split data into folds, and try each fold as validation and average the results



Useful for small datasets, but not used too frequently in deep learning

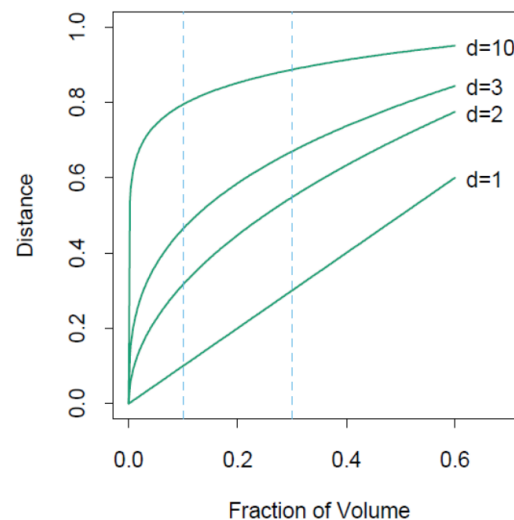
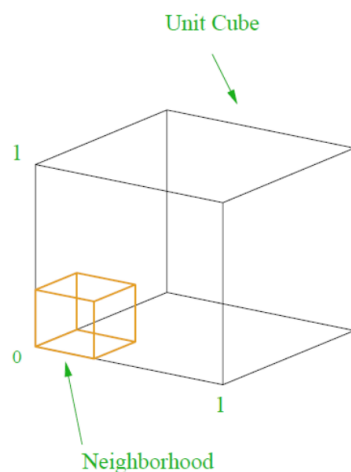
The Curse of Dimensionality

- Let's go back to data representation
- Pixels are actually never used...!
- This is due to the curse of dimensionality.
 - When the dimensionality increases, the volume of the space increases so fast that the available data become sparse



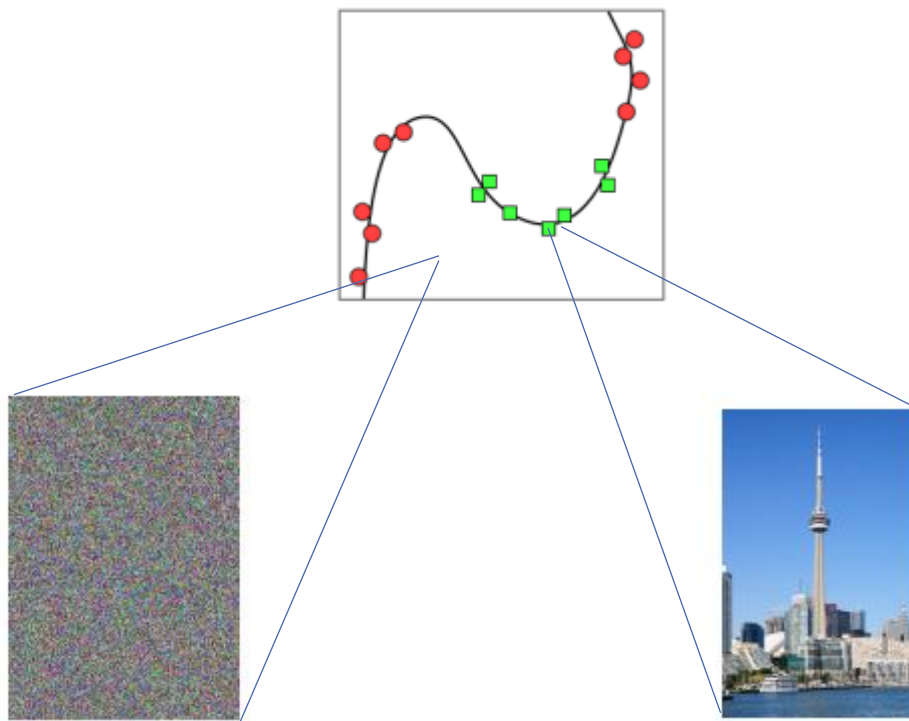
The Curse of Dimensionality

- In high dimensions, “most” points are far apart.
- If we want the nearest neighbor to be closer than ϵ , how many points do we need to guarantee it?
- The volume of a single ball radius ϵ is $\mathcal{O}(\epsilon^d)$
- The total volume of $[0, 1]^d$ is 1.
- Therefore, $\mathcal{O}\left(\left(\frac{1}{\epsilon}\right)^d\right)$ balls are needed to cover the volume.



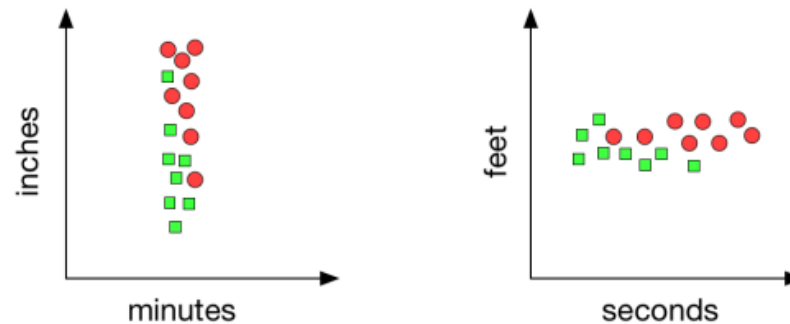
The Curse of Dimensionality

- Saving grace: some datasets (e.g., images) may have low intrinsic dimension (i.e., lie on or near a low-dimensional manifold).
- So nearest neighbors sometimes still works in high dimensions



Normalization

- Nearest Neighbors can be sensitive to the ranges of different features
- Often, the units are arbitrary:



- Simple fix: **normalize** each dimension to be zero mean and unit variance
- Caution: Depending on the problem the scale might be important!