# CMP784

## DEEP LEARNING

Lecture #2 – Machine Learning Overview

Puppy or bagel? // Karen Zack

Erkut Erdem // Hacettepe University // Fall 2023

HACETTEPE UNIVERSITY COMPUTER VISION LAB

# Previously on CMP784

- what is deep learning

- a brief history of deep learning

- compositionality

- end-to-end learning

- distributed representations

# Lecture overview

- what is learning?
- types of machine learning problems
- image classification
- linear regression
- generalization
- cross-validation
- maximum likelihood estimation

- **Disclaimer:** Much of the material and slides for this lecture were borrowed from
  - Bernhard Schölkopf's MLSS 2017 lecture,
  - Tommi Jaakkola's 6.867 class,
  - Fei-Fei Li, Andrej Karpathy and Justin Johnson's CS231n class
  - Hongsheng Li's ELEG5491 class
  - Justin Johnson's EECS 498/598 class

# Pre-quiz

- Due Date: 5pm, Sat, Oct 14, 2023.

- Each student enrolled to CMP784 must complete and pass this quiz!

# Good news, everyone!

- Practical 1 will be out next week
    - Multi-layer perceptrons
    - Word embeddings
    - Due November 1, 2023
    - Get familiar with Google CoLab

# What is learning?

# Two definitions of learning

- "Learning is the acquisition of knowledge about the world."
  Kupfermann (1985)


- "Learning is an adaptive change in behavior caused by experience."
  Shepherd (1988)

# Machine learning aims at learning a function

- Given an input value or vector, a function assigns it with a value or vector.

- "One-to-many" mapping is not a function. "Many-to-one" mapping is a function.

**Not a function**          **A function**

- Note that a function can have a vector output or matrix output. For instance, the following formula is still a function

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = f\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} x_1 + x_2 \\ x_1 x_2 \end{bmatrix}$$

# Function estimation

- We are interested in predicting $y$ from input $x$ and assume there exists a function that describes the relationship between $y$ and $x$, e.g., $y = f(x)$.

- If the function $f$ 's parametric form is fixed, prediction function $f$ can be parametrized by a parameter vector $\theta$.

- Estimating $\hat{f}$ from a training set $D = \{(x_1^{train}, y_1), (x_2^{train}, y_2), \cdots, (x_n^{train}, y_n)\}$.

- With a better design of the parametric form of the function, the learner could achieve better performance.

- This design process typical involves domain knowledge.

# Empirical Inference

- Drawing conclusions from empirical data (observations, measurements)

- Example 1: scientific inference

# Empirical Inference

- Drawing conclusions from empirical data (observations, measurements)

- Example 1: scientific inference

# Empirical Inference

- Drawing conclusions from empirical data (observations, measurements)

- Example 1: scientific inference

$$y = \sum_i a_i \, k(x, x_i) + b$$

Leibniz, Weyl, Chaitin
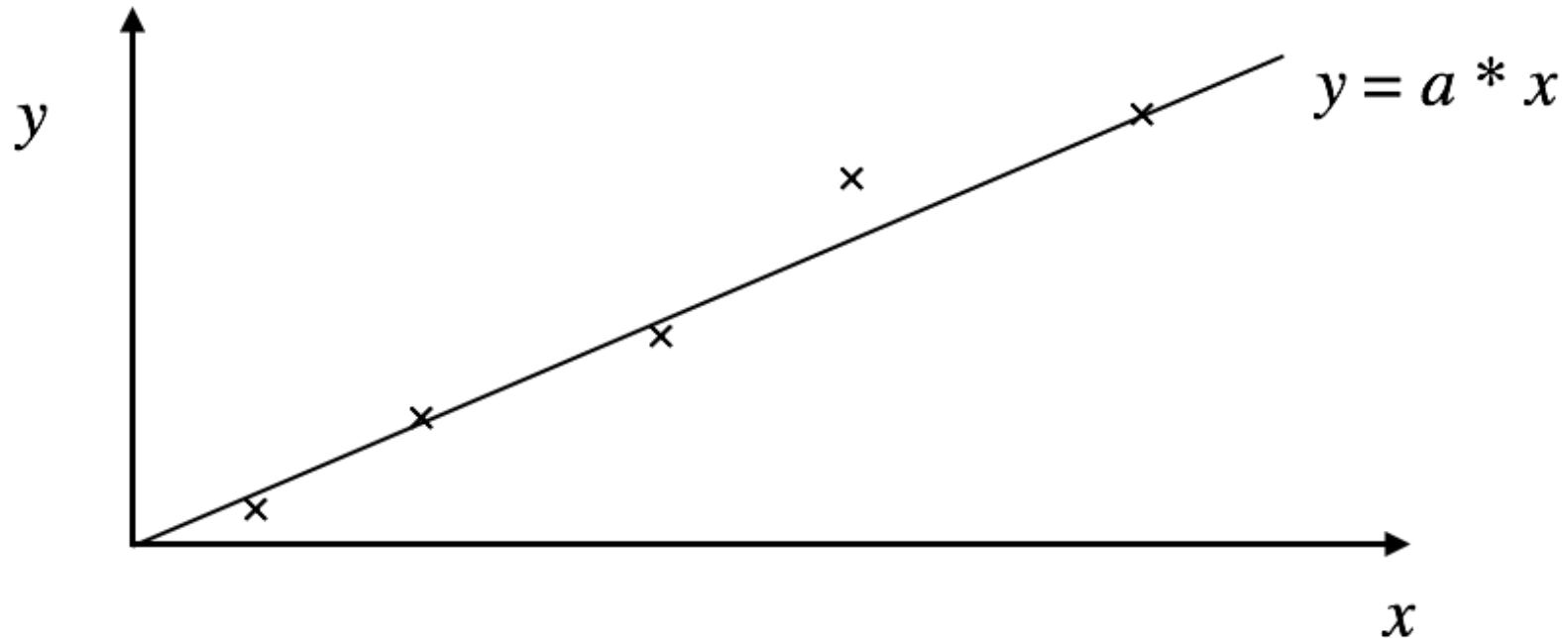
# Empirical Inference

- Drawing conclusions from empirical data (observations, measurements)

- Example 1: scientific inference

$$y = \sum_i a_i k(x, x_i) + b$$

Leibniz, Weyl, Chaitin

# Empirical Inference

- Example 2: perception

9

9

8

9

8

8

9

8

8

9

9

9

8

9

8

8

9

8

8

9

The choice of representation may determine whether the learning task is very easy or very difficult!

# Generalization

- *observe*
- What's next?

$$1, \quad 2, \quad 4, \quad 7, \ldots$$

$$+1 \quad +2 \quad +3$$

- $1,2,4,7,11,16,\ldots$: $a_{n+1} = a_n + n$ ("lazy caterer's sequence")
- $1,2,4,7,12,20,\ldots$: $a_{n+2} = a_{n+1} + a_n + 1$
- $1,2,4,7,13,24,\ldots$: "Tribonacci"-sequence
- $1,2,4,7,14,28$ : divisors of 28
- $1,2,4,7,1,1,5,\ldots$ : decimal expansions of $\pi$=3.14159… and
  $e$=2.718… interleaved *(thanks to O. Bousquet)*
- don't need $e$: 1247 appears at position 16992 in $\pi$
- The On-Line Encyclopedia of Integer Sequences: > 600 hits…

# Generalization, II

- Question: which continuation is correct ("generalizes")?

- *Answer?* There's no way to tell *("induction problem")*

- Question of statistical learning theory: how to come up with a law that generalizes *("demarcation problem")*

# Types of ML problems

# Types of machine learning problems

Based on the information available:

- Supervised learning

- Unsupervised learning

- Semi-supervised learning

- Reinforcement learning

# Supervised learning

- **Input:** $\{(\mathbf{x}, y)\}$

- **Task:** Predict target $y$ from input $\mathbf{x}$
  - **Classification:** Discrete output
  - **Regression:** Real-valued output



cat

Image classification



Face detection

# Unsupervised learning

- **Input:** $\{\mathbf{x}\}$

- **Task:** Reveal structure in the observed data
  - **Clustering:** Partition data into groups
  - **Feature extraction:** Learning meaningful features automatically
  - **Dimensionality reduction:** Learning a lower-dimensional representation of input

Agglomerative clustering

Topic modeling

Anomaly detection

# Semi-supervised learning

- **Input:**

    Few labeled examples $\{(\mathbf{x}, y)\}$

    Many unlabeled examples $\{\mathbf{x}\}$

- **Task:** Predict target $y$ from input $\mathbf{x}$
    - **Classification:** Discrete output
    - **Regression:** Real-valued output



interactive segmentation

Try to improve predictions based on examples by making use of the additional "unlabeled" examples

42

# Reinforcement learning

- **Input:**

  Interaction with an environment; the agent receives a numerical reward signal



action

Agent          Environment

observation, reward

- **Task:** A way of behaving that is very rewarding in the long run

- Goal is to estimate and maximize the long-term cumulative reward



TD-Gammon (Tesauro, 1990-1995)

# Types of machine learning problems



**How Much Information Does the Machine Need to Predict?**

Y LeCun

- **"Pure" Reinforcement Learning (cherry)**
  - ► The machine predicts a scalar reward given once in a while.
  - ► **A few bits for some samples**

- **Supervised Learning (icing)**
  - ► The machine predicts a category or a few numbers for each input
  - ► Predicting human-supplied data
  - ► **10→10,000 bits per sample**

- **Unsupervised/Predictive Learning (cake)**
  - ► The machine predicts any part of its input for any observed part.
  - ► Predicts future frames in videos
  - ► **Millions of bits per sample**

- **(Yes, I know, this picture is slightly offensive to RL folks. But I'll make it up)**

*"If intelligence was a cake, unsupervised learning would be the cake, supervised learning would be the icing on the cake, and reinforcement learning would be the cherry on the cake. We know how to make the icing and the cherry, but we don't know how to make the cake."*

– Yann LeCun
NeurIPS 2016 Keynote

44

# Image classification

— non-parametric vs. parametric models
— nearest neighbor classifier
— hyperparameter
— cross-validation

# Image Classification: a core task in Computer Vision



(assume given set of discrete labels)
{dog, cat, truck, plane, ...}

⟶ cat

# The problem: semantic gap

Images are represented as 3D arrays of numbers, with integers between [0, 255].

e.g.
300 x 100 x 3

(3 for 3 color channels RGB)



What the computer sees

# Challenges: Viewpoint Variation

# Challenges: Illumination

# Challenges: Deformation

# Challenges: Occlusion

# Challenges: Background clutter

# Challenges: Intraclass variation

# An image classifier

```python
def predict(image):
    # ????
    return class_label
```

Unlike e.g. sorting a list of numbers,

**no obvious way** to hard-code the algorithm for recognizing a cat, or other classes.

# Data-driven approach:
1. Collect a dataset of images and labels
2. Use Machine Learning to train an image classifier
3. Evaluate the classifier on a withheld set of test images

```python
def train(train_images, train_labels):
    # build a model for images -> labels...
    return model

def predict(model, test_images):
    # predict test_labels using the model...
    return test_labels
```

**Example training set**

# First classifier: Nearest Neighbor Classifier

```python
def train(train_images, train_labels):
    # build a model for images -> labels...
    return model


def predict(model, test_images):
    # predict test_labels using the model...
    return test_labels
```

Remember all training images and their labels

Predict the label of the most similar training image

Example dataset: **CIFAR-10**
**10** labels
**50,000** training images, each image is tiny: 32x32
**10,000** test images.

Example dataset: **CIFAR-10**
**10** labels
**50,000** training images
**10,000** test images.

For every test image (first column),
examples of nearest neighbors in rows

# How do we compare the images? What is the **distance metric**?

**L1 distance:** $$d_1(I_1, I_2) = \sum_p |I_1^p - I_2^p|$$

test image

| 56 | 32 | 10 | 18 |
|----|----|----|----|
| 90 | 23 | 128 | 133 |
| 24 | 26 | 178 | 200 |
| 2 | 0 | 255 | 220 |

\-

training image

| 10 | 20 | 24 | 17 |
|----|----|----|----|
| 8 | 10 | 89 | 100 |
| 12 | 16 | 178 | 170 |
| 4 | 32 | 233 | 112 |

=

pixel-wise absolute value differences

| 46 | 12 | 14 | 1 |
|----|----|----|----|
| 82 | 13 | 39 | 33 |
| 12 | 10 | 0 | 30 |
| 2 | 32 | 22 | 108 |

add → 456

# Nearest Neighbor classifier

```python
import numpy as np

class NearestNeighbor:
  def __init__(self):
    pass

  def train(self, X, y):
    """ X is N x D where each row is an example. Y is 1-dimension of size N """
    # the nearest neighbor classifier simply remembers all the training data
    self.Xtr = X
    self.ytr = y

  def predict(self, X):
    """ X is N x D where each row is an example we wish to predict label for """
    num_test = X.shape[0]
    # lets make sure that the output type matches the input type
    Ypred = np.zeros(num_test, dtype = self.ytr.dtype)

    # loop over all test rows
    for i in xrange(num_test):
      # find the nearest training image to the i'th test image
      # using the L1 distance (sum of absolute value differences)
      distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
      min_index = np.argmin(distances) # get the index with smallest distance
      Ypred[i] = self.ytr[min_index] # predict the label of the nearest example

    return Ypred
```

# Nearest Neighbor classifier

```python
import numpy as np

class NearestNeighbor:
  def __init__(self):
    pass

  def train(self, X, y):
    """ X is N x D where each row is an example. Y is 1-dimension of size N """
    # the nearest neighbor classifier simply remembers all the training data
    self.Xtr = X
    self.ytr = y

  def predict(self, X):
    """ X is N x D where each row is an example we wish to predict label for """
    num_test = X.shape[0]
    # lets make sure that the output type matches the input type
    Ypred = np.zeros(num_test, dtype = self.ytr.dtype)

    # loop over all test rows
    for i in xrange(num_test):
      # find the nearest training image to the i'th test image
      # using the L1 distance (sum of absolute value differences)
      distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
      min_index = np.argmin(distances) # get the index with smallest distance
      Ypred[i] = self.ytr[min_index] # predict the label of the nearest example

    return Ypred
```

remember the training data

# Nearest Neighbor classifier

```python
import numpy as np

class NearestNeighbor:
  def __init__(self):
    pass

  def train(self, X, y):
    """ X is N x D where each row is an example. Y is 1-dimension of size N """
    # the nearest neighbor classifier simply remembers all the training data
    self.Xtr = X
    self.ytr = y

  def predict(self, X):
    """ X is N x D where each row is an example we wish to predict label for """
    num_test = X.shape[0]
    # lets make sure that the output type matches the input type
    Ypred = np.zeros(num_test, dtype = self.ytr.dtype)

    # loop over all test rows
    for i in xrange(num_test):
      # find the nearest training image to the i'th test image
      # using the L1 distance (sum of absolute value differences)
      distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
      min_index = np.argmin(distances) # get the index with smallest distance
      Ypred[i] = self.ytr[min_index] # predict the label of the nearest example

    return Ypred
```

for every test image:
– find nearest train image
  with L1 distance
– predict the label of
  nearest training image

```python
import numpy as np

class NearestNeighbor:
  def __init__(self):
    pass

  def train(self, X, y):
    """ X is N x D where each row is an example. Y is 1-dimension of size N """
    # the nearest neighbor classifier simply remembers all the training data
    self.Xtr = X
    self.ytr = y

  def predict(self, X):
    """ X is N x D where each row is an example we wish to predict label for """
    num_test = X.shape[0]
    # lets make sure that the output type matches the input type
    Ypred = np.zeros(num_test, dtype = self.ytr.dtype)

    # loop over all test rows
    for i in xrange(num_test):
      # find the nearest training image to the i'th test image
      # using the L1 distance (sum of absolute value differences)
      distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
      min_index = np.argmin(distances) # get the index with smallest distance
      Ypred[i] = self.ytr[min_index] # predict the label of the nearest example

    return Ypred
```

Nearest Neighbor classifier

**Q: how does the classification speed depend on the size of the training data?**

```python
import numpy as np

class NearestNeighbor:
  def __init__(self):
    pass

  def train(self, X, y):
    """ X is N x D where each row is an example. Y is 1-dimension of size N """
    # the nearest neighbor classifier simply remembers all the training data
    self.Xtr = X
    self.ytr = y

  def predict(self, X):
    """ X is N x D where each row is an example we wish to predict label for """
    num_test = X.shape[0]
    # lets make sure that the output type matches the input type
    Ypred = np.zeros(num_test, dtype = self.ytr.dtype)

    # loop over all test rows
    for i in xrange(num_test):
      # find the nearest training image to the i'th test image
      # using the L1 distance (sum of absolute value differences)
      distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
      min_index = np.argmin(distances) # get the index with smallest distance
      Ypred[i] = self.ytr[min_index] # predict the label of the nearest example

    return Ypred
```

Nearest Neighbor classifier

Q: how does the classification speed depend on the size of the training data? **linearly :(**
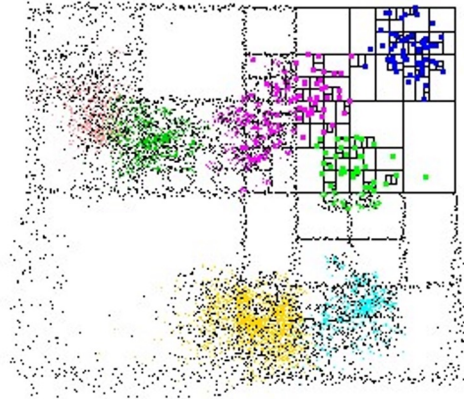
This is **backwards**:
- test time performance is usually much more important in practice.
- Deep Neural Networks flip this: expensive training, cheap test evaluation

# Aside: Approximate Nearest Neighbor
## find approximate nearest neighbors quickly



**ANN: A Library for Approximate Nearest Neighbor Searching**

**David M. Mount** and **Sunil Arya**

Version 1.1.2
Release Date: Jan 27, 2010

### What is ANN?

ANN is a library written in C++, which supports data structures and algorithms for both exact and approximate nearest neighbor searching in arbitrarily high dimensions.

In the nearest neighbor problem a set of data points in d-dimensional space is given. These points are preprocessed into a data structure, so that given any query point q, the nearest or generally k nearest points of P to q can be reported efficiently. The distance between two points can be defined in many ways. ANN assumes that distances are measured using any class of distance functions called Minkowski metrics. These include the well known Euclidean distance, Manhattan distance, and max distance.

Based on our own experience, ANN performs quite efficiently for point sets ranging in size from thousands to hundreds of thousands, and in dimensions as high as 20. (For applications in significantly higher dimensions, the results are rather spotty, but you might try it anyway.)

The library implements a number of different data structures, based on kd-trees and box-decomposition trees, and employs a couple of different search strategies.

The library also comes with test programs for measuring the quality of performance of ANN on any particular data sets, as well as programs for visualizing the structure of the geometric data structures.

---

**FLANN - Fast Library for Approximate Nearest Neighbors**

- Home
- News
- Publications
- Download
- Changelog
- Repository

#### What is FLANN?

FLANN is a library for performing fast approximate nearest neighbor searches in high dimensional spaces. It contains a collection of algorithms we found to work best for nearest neighbor search and a system for automatically choosing the best algorithm and optimum parameters depending on the dataset.

FLANN is written in C++ and contains bindings for the following languages: C, MATLAB and Python.

#### News

- (14 December 2012) Version 1.8.0 is out bringing incremental addition/removal of points to/from indexes
- (20 December 2011) Version 1.7.0 is out bringing two new index types and several other improvements.
- You can find binary installers for FLANN on the Point Cloud Library project page. Thanks to the PCL developers!
- Mac OS X users can install flann though MacPorts (thanks to Mark Moll for maintaining the Portfile)
- New release introducing an easier way to use custom distances, kd-tree implementation optimized for low dimensionality search and experimental MPI support
- New release introducing new C++ templated API, thread-safe search, save/load of indexes and more.
- The FLANN license was changed from LGPL to BSD.

#### How fast is it?

In our experiments we have found FLANN to be about one order of magnitude faster on many datasets (in query time), than previously available approximate nearest neighbor search software.

#### Publications

More information and experimental results can be found in the following papers:
- Marius Muja and David G. Lowe. **"Scalable Nearest Neighbor Algorithms for High Dimensional Data"**. Pattern Analysis and Machine Intelligence (PAMI), Vol. 36, 2014. [PDF] [BibTeX]
- Marius Muja and David G. Lowe. **"Fast Matching of Binary Features"**. Conference on Computer and Robot Vision (CRV) 2012. [PDF] [BibTeX]
- Marius Muja and David G. Lowe, **"Fast Approximate Nearest Neighbors with Automatic Algorithm Configuration"**, in International Conference on Computer Vision Theory and Applications (VISAPP'09), 2009 [PDF] [BibTeX]

# The choice of distance is a **hyperparameter** common choices:

L1 (Manhattan) distance

$$d_1(I_1, I_2) = \sum_p |I_1^p - I_2^p|$$

L2 (Euclidean) distance

$$d_2(I_1, I_2) = \sqrt{\sum_p (I_1^p - I_2^p)^2}$$

# k-Nearest Neighbor

find the k nearest images, have them vote on the label

the data          NN classifier          5-NN classifier



http://en.wikipedia.org/wiki/K-nearest_neighbors_algorithm

Example dataset: **CIFAR-10**
**10** labels
**50,000** training images
**10,000** test images.

For every test image (first column),
examples of nearest neighbors in rows

the data | NN classifier | 5-NN classifier

Q: what is the accuracy of the nearest neighbor classifier on the training data, when using the Euclidean distance?

the data | NN classifier | 5-NN classifier

Q2: what is the accuracy of the **k**-nearest neighbor classifier on the training data?

What is the best **distance** to use?
What is the best value of **k** to use?

i.e. how do we set the **hyperparameters**?

What is the best **distance** to use?
What is the best value of **k** to use?

i.e. how do we set the **hyperparameters**?

Very problem-dependent.
Must try them all out and see what works best.

Try out what hyperparameters work best on test set.

Trying out what hyperparameters work best on test set:
Very bad idea. The test set is a proxy for the generalization performance!
Use only **VERY SPARINGLY**, at the end.

| train data | test data |

**Validation data**
use to tune hyperparameters

**Cross-validation**

cycle through the choice of which fold is the validation fold, average results.

Cross-validation on k

Example of
5-fold cross-validation
for the value of **k**.

Each point: single
outcome.

The line goes
through the mean, bars
indicated standard
deviation

(Seems that k ~= 7 works
the best for this data)

# k-Nearest Neighbor on images **never used.**

- terrible performance at test time
- distance metrics on level of whole images can be very unintuitive



| original | shifted | messed up | darkened |

(all 3 images have same L2 distance to the one on the left)

# Nearest Neighbor with ConvNet features works well!



Devlin et al, "Exploring Nearest Neighbor Approaches for Image Captioning", 2015

# Nearest Neighbor with ConvNet features works well!

Example: Image Captioning with Nearest Neighbor

A bedroom with a bed and a couch.

A cat sitting in a bathroom sink.

A train is stopped at a train station.

A wooden bench in front of a building.

Devlin et al, "Exploring Nearest Neighbor Approaches for Image Captioning", 2015

# The learning problem

— linear classification
— hypothesis class, estimation algorithm
— loss and estimation criterion
— sampling, empirical and expected losses

# The Learning Problem

**Digit Recognition**



**Image Classification**



cat

car

frog

- Steps
  - —entertain a (biased) set of possibilities (hypothesis class)
  - —adjust predictions based on available examples (estimation)
  - —rethink the set of possibilities (model selection)

- Principles of learning are "universal"
  - —society (e.g., scientific community)
  - —animal (e.g., human)
  - —machine

# Hypothesis class

- Representation: examples are binary vecto[...]th $d = 64$

$$\mathbf{x} = [111\dots0001]^T$$

and labels $y \in \{-1, 1\}$ ("no","yes")

- The mapping from examples to labels is a "**linear classifier**"

$$\hat{y} = \text{sign}(\theta \cdot \mathbf{x}) = \text{sign}(\theta_1 x_1 + \dots + \theta_d x_d)$$

where $\theta$ is a vector of **parameters** we have to learn from examples.

# Linear classifier/experts

- We can understand the simple linear classifier

$$\hat{y} = \mathrm{sign}\,(\,\theta \cdot \mathbf{x}\,) = \mathrm{sign}\,(\,\theta_1 x_1 + \ldots + \theta_d x_d\,)$$

as a way of combining expert opinion (in this case simple binary features)

majority rule

$$\hat{y} = \mathrm{sign}\,(\,\theta_1 x_1 + \ldots + \theta_d x_d\,)$$

"votes"

combined "votes"

$\theta_1$  $\theta_2$  $\theta_d$

Expert 1

$\cdots$

$x_1$  $x_2$  $x_d$

1  0  1

# Estimation

| $\mathbf{x}$ | $y$ |
|---|---|
| 0111111100111001000000010000000100111111101110111110011101111110001 | +1 |
| 0001111100000001100000111000000110011111110111111100111111110000011 | +1 |
| 1111111000000110000011000111111100000001111000001111100011011111111 | −1 |
| . . . | . . . |

- How do we adjust the parameters $\theta$ based on the labeled examples?

$$\hat{y} = \text{sign}\left(\,\theta \cdot \mathbf{x}\,\right)$$

For example, we can simply refine/update the parameters whenever we make a mistake (**perceptron algorithm**):

$$\theta_i \leftarrow \theta_i + y\,x_i,\ i = 1,\dots,d \quad \text{if prediction was wrong}$$

# Evaluation

- Does the simple mistake driven algorithm work?



(average classification error as a function of the number of examples and labels seen so far)

# Illustration of Convergence

• Convergence of the perceptron learning algorithm

# Linear classifier: image classification



image    parameters

$$f(x, W)$$

**10** numbers, indicating class scores

**[32x32x3]**
array of numbers 0...1
(3072 numbers total)

# Linear classifier: image classification

$$f(x, W) = Wx$$



**[32x32x3]**
array of numbers 0…1

**10** numbers,
indicating class
scores

# Linear classifier: image classification

$$f(x, W) = Wx$$

**3072x1**

**10x1**   **10x3072**

**[32x32x3]**
array of numbers 0...1

**10** numbers,
indicating class
scores

parameters, or "weights"

# Linear classifier: image classification

$$f(x, W) = Wx \quad (+b)$$

**10x1**     **10x3072**     **3072x1**     **10x1**



**[32x32x3]**
array of numbers 0...1

**10** numbers, indicating class scores

parameters, or "weights"

# Example with an image with 4 pixels, and 3 classes (cat/dog/ship)

stretch pixels into single column



| | | | |
|---|---|---|---|
| 0.2 | -0.5 | 0.1 | 2.0 |
| 1.5 | 1.3 | 2.1 | 0.0 |
| 0 | 0.25 | 0.2 | -0.3 |

$W$

input image

| |
|---|
| 56 |
| 231 |
| 24 |
| 2 |

$x_i$

$+$

| |
|---|
| 1.1 |
| 3.2 |
| -1.2 |

$b$

$\longrightarrow$

| | |
|---|---|
| -96.8 | cat score |
| 437.9 | dog score |
| 61.95 | ship score |

$f(x_i; W, b)$

# Interpreting a Linear Classifier (Visual Viewpoint)



$$f(x_i, W, b) = Wx_i + b$$

Q: what does the linear classifier do, in English?

# Interpreting a Linear Classifier (Visual Viewpoint)



$$f(x_i, W, b) = Wx_i + b$$

Example trained weights of a linear classifier trained on CIFAR-10:

# Interpreting a Linear Classifier (Visual Viewpoint)



Linear classifier has one "template" per category

A single template cannot capture multiple modes of the data

e.g. horse template has 2 heads!

# Interpreting a Linear Classifier (Geometric Viewpoint)

$$f(x_i, W, b) = Wx_i + b$$



airplane classifier

car classifier

deer classifier

**[32x32x3]**
array of numbers 0...1
(3072 numbers total)

# Interpreting a Linear Classifier (Geometric Viewpoint)



$$f(x_i, W, b) = W x_i + b$$

car score increases this way

car classifier

airplane score increases this way

airplane classifier

deer classifier

deer score increases this way

**[32x32x3]**

array of numbers 0...1
(3072 numbers total)

# Model selection

- The simple linear classifier cannot solve all the problems

**Class 1**:
First and third quadrants

**Class 2**:
Second and fourth quadrants

**Class 1**:
1 <= L2 norm <= 2

**Class 2**:
Everything else

**Class 1**:
Three modes

**Class 2**:
Everything else



- Can we rethink the approach to do even better?
- We can, for example, add "polynomial experts"

$$\hat{y} = \text{sign} \left( \theta_1 x_1 + \ldots + \theta_d x_d + \textcolor{red}{\theta_{12} x_1 x_2} + \ldots \right)$$

# Model selection (cont'd)



linear

2<sup>nd</sup> order polynomial

4<sup>th</sup> order polynomial

8<sup>th</sup> order polynomial

# Review: The learning problem

**Image Cla**



cat

car

frog

- **Hypothesis class:** we consider some **restricted** set $\mathcal{F}$ of mappings $f : \mathcal{X} \to \mathcal{L}$ from images to labels

- **Estimation:** on the basis of a training set of examples and labels, $\{(\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_n, y_n)\}$, we find an estimate $\hat{f} \in \mathcal{F}$

- **Evaluation:** we measure how well $\hat{f}$ **generalizes** to yet unseen examples, i.e., whether $\hat{f}(\mathbf{x}_{new})$ agrees with $y_{new}$

# Hypothesis and estimation

- We used a simple linear classifier, a parameterized mapping $f(\mathbf{x}; \theta)$ from images $\mathcal{X}$ to labels $\mathcal{L}$, to solve a binary image classification problem (2's vs 3's):

$$\hat{y} \;=\; f(\mathbf{x}; \theta) = \mathsf{sign}\big(\,\theta \cdot \mathbf{x}\,\big)$$

  where $\mathbf{x}$ is a pixel image and $\hat{y} \in \{-1, 1\}$.

- The parameters $\theta$ were adjusted on the basis of the training examples and labels according to a simple mistake driven update rule (written here in a vector form)

$$\theta \leftarrow \theta + y_i \mathbf{x}_i, \quad \text{whenever} \quad y_i \neq \mathsf{sign}\big(\,\theta \cdot \mathbf{x}_i\,\big)$$

- The update rule attempts to minimize the number of errors that the classifier makes on the training examples

# Estimation criterion

- We can formulate the binary classification problem more explicitly by defining a **zero-one loss:**

$$\mathsf{Loss}(y, \hat{y}) = \begin{cases} 0, y = \hat{y} \\ 1, y \neq \hat{y} \end{cases}$$

so that

$$\frac{1}{n}\sum_{i=1}^{n}\mathsf{Loss}(y_i, \hat{y}_i) = \frac{1}{n}\sum_{i=1}^{n}\mathsf{Loss}(y_i, f(\mathbf{x}_i; \theta))$$

gives the fraction of prediction errors on the training set.

- This is a function of the parameters $\theta$ and we can try to minimize it directly.

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



| | | | |
|---|---|---|---|
| cat | **3.2** | 1.3 | 2.2 |
| car | 5.1 | **4.9** | 2.5 |
| frog | -1.7 | 2.0 | **-3.1** |

**Multiclass SVM loss:**

Given an example $(x_i, y_i)$
where $x_i$ is the image and
where $y_i$ is the (integer) label,

and using the shorthand for the
scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



|  | | | |
|---|---|---|---|
| cat | **3.2** | 1.3 | 2.2 |
| car | 5.1 | **4.9** | 2.5 |
| frog | -1.7 | 2.0 | **-3.1** |
| Losses: | 2.9 | | |

**Multiclass SVM loss:**

Given an example $(x_i, y_i)$
where $x_i$ is the image and
where $y_i$ is the (integer) label,

and using the shorthand for the
scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

= max(0, 5.1 - 3.2 + 1)
    +max(0, -1.7 - 3.2 + 1)
= max(0, 2.9) + max(0, -3.9)
= 2.9 + 0
= 2.9

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:

**Multiclass SVM loss:**

Given an example $(x_i, y_i)$
where $x_i$ is the image and
where $y_i$ is the (integer) label,

and using the shorthand for the
scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

|  | | | |
|---|---|---|---|
| cat | **3.2** | 1.3 | 2.2 |
| car | 5.1 | **4.9** | 2.5 |
| frog | -1.7 | 2.0 | **-3.1** |
| Losses: | 2.9 | 0 | |

= max(0, 1.3 - 4.9 + 1)
  +max(0, 2.0 - 4.9 + 1)
= max(0, -2.6) + max(0, -1.9)
= 0 + 0
= 0

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



| | cat | car | frog |
|---|---|---|---|
| cat | **3.2** | 1.3 | 2.2 |
| car | 5.1 | **4.9** | 2.5 |
| frog | -1.7 | 2.0 | **-3.1** |
| Losses: | 2.9 | 0 | 10.9 |

**Multiclass SVM loss:**

Given an example $(x_i, y_i)$
where $x_i$ is the image and
where $y_i$ is the (integer) label,

and using the shorthand for the
scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

= max(0, 2.2 - (-3.1) + 1)
  +max(0, 2.5 - (-3.1) + 1)
= max(0, 5.3) + max(0, 5.6)
= 5.3 + 5.6
= 10.9

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



| | cat | car | frog |
|---|---|---|---|
| cat | **3.2** | 1.3 | 2.2 |
| car | 5.1 | **4.9** | 2.5 |
| frog | -1.7 | 2.0 | **-3.1** |
| Losses: | 2.9 | 0 | 10.9 |

**Multiclass SVM loss:**

Given an example $(x_i, y_i)$ where $x_i$ is the image and where $y_i$ is the (integer) label,

and using the shorthand for the scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$

and the full training loss is the mean over all examples in the training data:

$$L = \frac{1}{N} \sum_{i=1}^{N} L_i$$

L = (2.9 + 0 + 10.9)/3
= 4.6

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



|  | cat | car | frog |
|------|------|------|------|
| cat | **3.2** | 1.3 | 2.2 |
| car | 5.1 | **4.9** | 2.5 |
| frog | -1.7 | 2.0 | **-3.1** |
| Losses: | 2.9 | 0 | 10.9 |

**Multiclass SVM loss:**

Given an example $(x_i, y_i)$ where $x_i$ is the image and where $y_i$ is the (integer) label,

and using the shorthand for the scores vector: $s = f(x_i, W)$

the SVM loss has the form:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

Q: what is the min/max possible loss?

There is something missing in the loss:

$$f(x, W) = Wx$$

$$L = \frac{1}{N} \sum_{i=1}^{N} \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1)$$

e.g. suppose that we found a W such that L = 0.
Is this W unique?

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$



| | cat | car | frog |
|---|---|---|---|
| cat | **3.2** | 1.3 | 2.2 |
| car | 5.1 | **4.9** | 2.5 |
| frog | -1.7 | 2.0 | **-3.1** |
| Losses: | 2.9 | 0 | |

Before:
= max(0, 1.3 - 4.9 + 1)
  +max(0, 2.0 - 4.9 + 1)
= max(0, -2.6) + max(0, -1.9)
= 0 + 0
= 0

With W twice as large:
= max(0, 2.6 - 9.8 + 1)
  +max(0, 4.0 - 9.8 + 1)
= max(0, -6.2) + max(0, -4.8)
= 0 + 0
= 0

117

# Weight Regularization

$$L = \frac{1}{N} \sum_{i=1}^{N} \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1) + \boxed{\lambda R(W)}$$

In common use:

**L2 regularization**

L1 regularization

Elastic net (L1 + L2)

Max norm regularization

Dropout (will see later)

Batch normalization (will see later)

$$R(W) = \sum_k \sum_l W_{k,l}^2$$

$$R(W) = \sum_k \sum_l |W_{k,l}|$$

$$R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$$

# L2 regularization: motivation

$$x = [1, 1, 1, 1]$$

$$w_1 = [1, 0, 0, 0]$$
$$w_2 = [0.25, 0.25, 0.25, 0.25]$$

$$w_1^T x = w_2^T x = 1$$

# Estimation criterion (revisited)

- We have reduced the estimation problem to a minimization problem

$$\text{find } \theta \text{ that minimizes} \quad \frac{1}{n} \sum_{i=1}^{n} \overbrace{\text{Loss}\big(y_i, f(\mathbf{x}_i; \theta)\big)}^{\text{empirical loss}}$$

  - valid for any parameterized class of mappings from examples to predictions
  - valid when the predictions are discrete labels, real valued, or other provided that the loss is defined appropriately
  - may be ill-posed (under-constrained) as stated

- But why is it sensible to minimize the empirical loss in the first place since we are only interested in the performance on new examples?

# Training and test performance: sampling

- We assume that each training **and** test example-label pair, $(\mathbf{x}, y)$ is drawn **independently at random** from the **same** but unknown population of examples and labels.

- We can represent this population as a joint probability distribution $P(\mathbf{x}, y)$ so that each training/test example is a **sample** from this distribution

$$(\mathbf{x}_i, y_i) \sim P$$

Indyk

Barzilay

Collins

Jaakkola

cat

car

frog

Indy

# Training and test performance: sampling

- We assume that each training **and** test example-label pair, $(\mathbf{x}, y)$ is drawn **independently at random** from the **same** but unknown population of examples and labels.

- We can represent this population as a joint probability distribution $P(\mathbf{x}, y)$ so that each training/test example is a **sample** from this distribution

$$(\mathbf{x}_i, y_i) \sim P$$

Indyk
Barzilay
Collins
Jaakkola

$$\text{Empirical (training) loss} \;\; = \;\; \frac{1}{n}\sum_{i=1}^{n} \text{Loss}\big(y_i, f(\mathbf{x}_i; \theta)\big)$$

$$\text{Expected (test) loss} \;\; = \;\; E_{(\mathbf{x}, y) \sim P}\big\{\text{Loss}\big(y, f(\mathbf{x}; \theta)\big)\big\}$$

- The training loss based on a few sampled examples and labels serves as a proxy for the test performance measured over the whole population

# Regression, example
# Linear regression
— Estimation, errors, analysis

# Regression

- The goal is to make quantitative (real valued) predictions on the basis of a (vector of) features or attributes

- Example: predicting vehicle fuel efficiency (mpg) from 8 attributes

| y | $\mathbf{x}$ | | | | |
|---|---|---|---|---|---|
| | cyls | disp | hp | weight | . . . |
| 18.0 | 8 | 307.0 | 130.00 | 3504 | . . . |
| 26.0 | 4 | 97.00 | 46.00 | 1835 | . . . |
| 33.5 | 4 | 98.00 | 83.00 | 2075 | . . . |
| . . . | | | | | |

- We need to
  - specify the class of functions (e.g., linear)
  - select how to measure prediction loss
  - solve the resulting minimization problem

# Linear regression



- We begin by considering linear regression (easy to extend to more complex predictions later on)

$$f : \mathcal{R} \to \mathcal{R} \qquad f(x; \mathbf{w}) = w_0 + w_1 x$$

$$f : \mathcal{R}^d \to \mathcal{R} \qquad f(\mathbf{x}; \mathbf{w}) = w_0 + w_1 x_1 + \ldots w_d x_d$$

where $\mathbf{w}$ are **parameters** we need to set.

# Linear reg u

$$f : \mathcal{R} \to \mathcal{R} \qquad f(x; \mathbf{w}) = w_0 + w_1 x$$

$$f : \mathcal{R}^d \to \mathcal{R} \qquad f(\mathbf{x}; \mathbf{w}) = w_0 + w_1 x_1 + \ldots w_d x_d$$

- We can measure the prediction loss in terms of squared error, $\text{Loss}(y, \hat{y}) = (y - \hat{y})^2$, so that the empirical loss on $n$ training samples becomes mean squared error

$$J_n(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^{n} \big( y_i - f(\mathbf{x}_i; \mathbf{w}) \big)^2$$

# Linear regression: estimation

- We have to minimize the **empirical** squared loss

$$J_n(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^{n} \left( y_i - f(\mathbf{x}_i; \mathbf{w}) \right)^2$$

$$= \frac{1}{n} \sum_{i=1}^{n} \left( y_i - w_0 - w_1 x_i \right)^2 \quad \text{(1-dim)}$$

- By setting the derivatives with respect to $w_1$ and $w_0$ to zero, we get necessary conditions for the "optimal" parameter values

$$\frac{\partial}{\partial w_1} J_n(\mathbf{w}) = 0$$

$$\frac{\partial}{\partial w_0} J_n(\mathbf{w}) = 0$$

# Optimality conditions: derivation

$$\frac{\partial}{\partial w_1} J_n(\mathbf{w}) \;=\; \frac{\partial}{\partial w_1} \frac{1}{n} \sum_{i=1}^{n} (y_i - w_0 - w_1 x_i)^2$$

# Optimality conditions: derivation

$$\frac{\partial}{\partial w_1} J_n(\mathbf{w}) = \frac{\partial}{\partial w_1} \frac{1}{n} \sum_{i=1}^{n} (y_i - w_0 - w_1 x_i)^2$$

$$= \frac{1}{n} \sum_{i=1}^{n} \frac{\partial}{\partial w_1} (y_i - w_0 - w_1 x_i)^2$$

# Optimality conditions: derivation

$$
\begin{aligned}
\frac{\partial}{\partial w_1} J_n(\mathbf{w}) \;\; &= \;\; \frac{\partial}{\partial w_1} \frac{1}{n} \sum_{i=1}^{n} \left( y_i - w_0 - w_1 x_i \right)^2 \\[2ex]
&= \;\; \frac{1}{n} \sum_{i=1}^{n} \frac{\partial}{\partial w_1} \left( y_i - w_0 - w_1 x_i \right)^2 \\[2ex]
&= \;\; \frac{2}{n} \sum_{i=1}^{n} \left( y_i - w_0 - w_1 x_i \right) \frac{\partial}{\partial w_1} (y_i - w_0 - w_1 x_i)
\end{aligned}
$$

# Optimality conditions: derivation

$$
\begin{aligned}
\frac{\partial}{\partial w_1} J_n(\mathbf{w}) &= \frac{\partial}{\partial w_1} \frac{1}{n} \sum_{i=1}^{n} (y_i - w_0 - w_1 x_i)^2 \\
&= \frac{1}{n} \sum_{i=1}^{n} \frac{\partial}{\partial w_1} (y_i - w_0 - w_1 x_i)^2 \\
&= \frac{2}{n} \sum_{i=1}^{n} (y_i - w_0 - w_1 x_i) \frac{\partial}{\partial w_1} (y_i - w_0 - w_1 x_i) \\
&= \frac{2}{n} \sum_{i=1}^{n} (y_i - w_0 - w_1 x_i)(-x_i) = 0
\end{aligned}
$$

# Optimality conditions: derivation

$$
\begin{aligned}
\frac{\partial}{\partial w_1} J_n(\mathbf{w}) &= \frac{\partial}{\partial w_1} \frac{1}{n} \sum_{i=1}^{n} (y_i - w_0 - w_1 x_i)^2 \\
&= \frac{1}{n} \sum_{i=1}^{n} \frac{\partial}{\partial w_1} (y_i - w_0 - w_1 x_i)^2 \\
&= \frac{2}{n} \sum_{i=1}^{n} (y_i - w_0 - w_1 x_i) \frac{\partial}{\partial w_1} (y_i - w_0 - w_1 x_i) \\
&= \frac{2}{n} \sum_{i=1}^{n} (y_i - w_0 - w_1 x_i)(-x_i) = 0 \\
\frac{\partial}{\partial w_0} J_n(\mathbf{w}) &= \frac{2}{n} \sum_{i=1}^{n} (y_i - w_0 - w_1 x_i)(-1) = 0
\end{aligned}
$$

# Linear regression: matrix notation

- We can express the solution a bit more generally by resorting to a matrix notation

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \cdots \\ y_n \end{bmatrix}, \ \mathbf{X} = \begin{bmatrix} 1 & x_1 \\ \cdots & \cdots \\ 1 & x_n \end{bmatrix}, \ \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

so that

$$\frac{1}{n}\sum_{t=1}^{n}(y_t - w_0 - w_1 x_t)^2 \ = \ \frac{1}{n}\left\Vert \begin{bmatrix} y_1 \\ \cdots \\ y_n \end{bmatrix} - \begin{bmatrix} 1 & x_1 \\ \cdots & \cdots \\ 1 & x_n \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} \right\Vert^2$$

$$= \ \frac{1}{n}\Vert \mathbf{y} - \mathbf{X}\mathbf{w} \Vert^2$$

# Linear regression: solution

- By setting the derivatives of $\|\mathbf{y} - \mathbf{Xw}\|^2/n$ to zero, we get the same optimality conditions as before, now expressed in a matrix form

$$\frac{\partial}{\partial \mathbf{w}} \frac{1}{n} \|\mathbf{y} - \mathbf{Xw}\|^2 \ = \ \frac{\partial}{\partial \mathbf{w}} \frac{1}{n} (\mathbf{y} - \mathbf{Xw})^T (\mathbf{y} - \mathbf{Xw})$$

# Linear regression: solution

- By setting the derivatives of $\|\mathbf{y} - \mathbf{Xw}\|^2/n$ to zero, we get the same optimality conditions as before, now expressed in a matrix form

$$\frac{\partial}{\partial \mathbf{w}} \frac{1}{n} \|\mathbf{y} - \mathbf{Xw}\|^2 = \frac{\partial}{\partial \mathbf{w}} \frac{1}{n} (\mathbf{y} - \mathbf{Xw})^T (\mathbf{y} - \mathbf{Xw})$$

$$= \frac{2}{n} \mathbf{X}^T (\mathbf{y} - \mathbf{Xw})$$

# Linear regression: solution

- By setting the derivatives of $\|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2/n$ to zero, we get the same optimality conditions as before, now expressed in a matrix form

$$
\begin{aligned}
\frac{\partial}{\partial \mathbf{w}} \frac{1}{n} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 &= \frac{\partial}{\partial \mathbf{w}} \frac{1}{n}(\mathbf{y} - \mathbf{X}\mathbf{w})^T(\mathbf{y} - \mathbf{X}\mathbf{w}) \\
&= \frac{2}{n}\mathbf{X}^T(\mathbf{y} - \mathbf{X}\mathbf{w}) \\
&= \frac{2}{n}(\mathbf{X}^T\mathbf{y} - \mathbf{X}^T\mathbf{X}\mathbf{w}) = \mathbf{0}
\end{aligned}
$$

which gives

$$
\hat{\mathbf{w}} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}
$$

- The solution is a linear function of the outputs $y$

# Alternative: Gradient Descent Algorithm

- $\Delta = -\frac{1}{T} \sum_t \nabla_{\boldsymbol{\theta}} l(f(\mathbf{x}^{(t)}; \boldsymbol{\theta}), y^{(t)}) - \lambda \nabla_{\boldsymbol{\theta}} \Omega(\boldsymbol{\theta})$

- One straightforward method: **gradient descent**

  - $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \Delta$ (initialize $\boldsymbol{\theta}$ e.g. randomly)

  - repeatedly update $\boldsymbol{\theta}$ based on the gradient

$$\Delta = -\frac{1}{T} \sum_t \nabla_{\boldsymbol{\theta}} l(f(\mathbf{x}^{(t)}; \boldsymbol{\theta}), y^{(t)}) - \lambda \nabla_{\boldsymbol{\theta}} \Omega(\boldsymbol{\theta})$$
$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \Delta$$

- $\alpha$ is the **learning rate**

5

# Effect of learning rate $\lambda$



- Large $\lambda$ => Fast convergence but larger residual error
  Also possible oscillations

- Small $\lambda$ => Slow convergence but small residual error

# Local and Global Optima

# Stochastic Gradient Descent

$$l(f(\mathbf{x}^{(t)}; \boldsymbol{\theta}), y^{(t)})$$

$$\Omega(\boldsymbol{\theta})$$

- Two ways to generalize this for all examples in training set:

$$\Delta = -\frac{1}{T} \sum_t \nabla_{\boldsymbol{\theta}} l(f(\mathbf{x}^{(t)}; \boldsymbol{\theta}), y^{(t)}) - \lambda \nabla_{\boldsymbol{\theta}} \Omega(\boldsymbol{\theta})$$

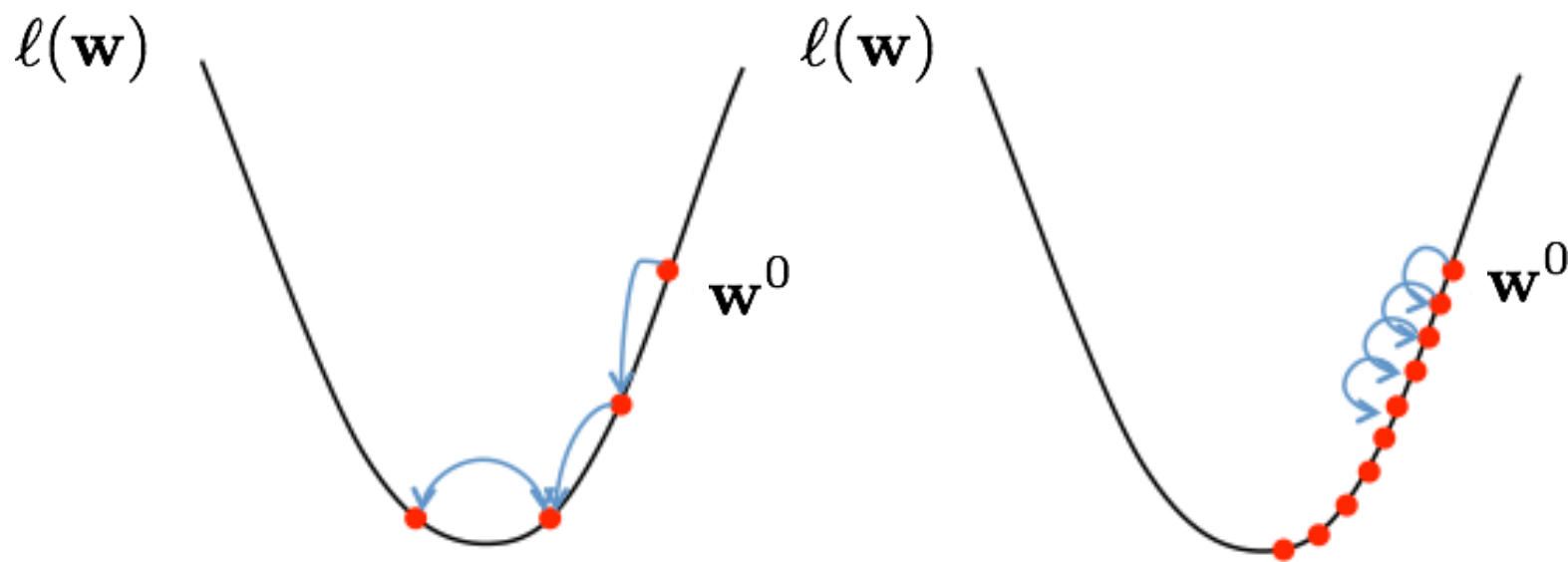1. **Batch updates:** sum or average updates across every example n, then change the parameter values

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \Delta$$

2. **Stochastic/online updates:** update the parameters for each training case in turn, according to its own gradients

$$\Delta = -\nabla_{\boldsymbol{\theta}} l(f(\mathbf{x}^{(t)}; \boldsymbol{\theta}), y^{(t)}) - \lambda \nabla_{\boldsymbol{\theta}} \Omega(\boldsymbol{\theta})$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \, \Delta$$

# Linear regression: generalization

- As the number of training examples increases our solution gets "better"



We'd like to understand the error a bit better

# Linear regression: types of errors

- **Structural error** measures the error introduced by the limited function class (infinite training data):

$$\min_{w_1, w_0} E_{(x,y) \sim P} \left( y - w_0 - w_1 x \right)^2 = E_{(x,y) \sim P} \left( y - w_0^* - w_1^* x \right)^2$$

where $(w_0^*, w_1^*)$ are the optimal linear regression parameters.

- **Approximation error** measures how close we can get to the optimal linear predictions with limited training data:

$$E_{(x,y) \sim P} \left( w_0^* + w_1^* x - \hat{w}_0 - \hat{w}_1 x \right)^2$$

where $(\hat{w}_0, \hat{w}_1)$ are the parameter estimates based on a small training set (therefore themselves random variables).

# Linear regression: error decomposition

- The expected error of our linear regression function decomposes into the sum of structural and approximation errors

$$E_{(x,y)\sim P}\left(y - \hat{w}_0 - \hat{w}_1 x\right)^2 =$$

$$E_{(x,y)\sim P}\left(y - w_0^* - w_1^* x\right)^2 +$$

$$E_{(x,y)\sim P}\left(w_0^* + w_1^* x - \hat{w}_0 - \hat{w}_1 x\right)^2$$

# Bias-Variance Tradeoff

- **Variance** of trained model: does it vary a lot if the training set changes

- **Bias** of trained model: is the average model close to the true solution?

- Generalization error can be seen as the sum of bias and the variance

# Parametric vs. non-parametric models

- **Parametric model:** its capacity is fixed and does not increase with the amount of training data
  - examples: linear classifier, neural network with fixed number of hidden units, etc.

- **Non-parametric model:** the capacity increases with the amount of training data
  - examples: k nearest neighbors classifier, neural network with adaptable hidden layer size, etc.

# Beyond linear regression models

— additive regression models, examples
— generalization and cross-validation
— population minimizer

# Linear regression

- Linear regression functions,

$$f : \mathcal{R} \to \mathcal{R} \qquad f(x; \mathbf{w}) = w_0 + w_1 x, \quad \text{or}$$

$$f : \mathcal{R}^d \to \mathcal{R} \qquad f(\mathbf{x}; \mathbf{w}) = w_0 + w_1 x_1 + \ldots + w_d x_d$$

  combined with the squared loss, are convenient because they are **linear in the parameters.**

—we get closed form estimates of the parameters

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

  where, for example, $\mathbf{y} = [y_1, \ldots, y_n]^T$.

—the resulting prediction errors $\epsilon_i = y_i - f(\mathbf{x}_i; \hat{\mathbf{w}})$ are uncorrelated with any linear function of the inputs $\mathbf{x}$.

—we can easily extend these to non-linear functions of the inputs while still keeping them linear in the parameters

# Beyond linear regression

- Example extension: $m^{th}$ order polynomial regression where $f : \mathcal{R} \rightarrow \mathcal{R}$ is given by

$$f(x; \mathbf{w}) = w_0 + w_1 x + \ldots + w_{m-1} x^{m-1} + w_m x^m$$

- — linear in the parameters, non-linear in the inputs
- — solution as before

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

where

$$\hat{\mathbf{w}} = \begin{bmatrix} \hat{w}_0 \\ \hat{w}_1 \\ \ldots \\ \hat{w}_m \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} 1 & x_1 & x_1^2 & \ldots & x_1^m \\ 1 & x_2 & x_2^2 & \ldots & x_2^m \\ \ldots & \ldots & \ldots & \ldots & \ldots \\ 1 & x_n & x_n^2 & \ldots & x_n^m \end{bmatrix}$$

# Polynomial regression



degree = 1

degree = 3

degree = 5

degree = 7

# Underfitting and Overfitting

Underfitting

Appropriate capacity

Overfitting

# Underfitting

Underfitting



$y$

$x_0$

- The learner cannot find a solution that fits training examples well
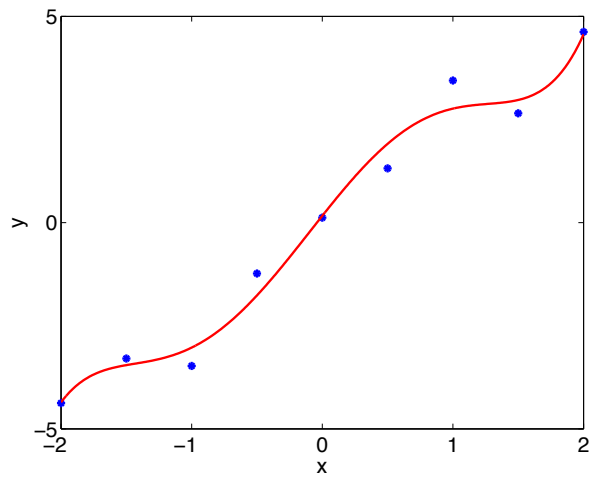  - For example, use linear regression to fit training examples $\{x^{(i)}, y^{(i)}\}$ where $y^{(i)}$ is an quadratic function of $x^{(i)}$.

- Underfitting means that the learner cannot capture some important aspects of the data.

- Reasons for underfitting happens
  - Model is not rich enough
  - Difficult to find the global optimum of the objective function on the training set or easy to get stuck at local minimum
  - Limitation on the computation resources (not enough training iterations of an iterative optimization procedure)

- Underfitting commonly happens in non-deep learning approaches with large scale training data and could be even a more serious problem than overfitting in some cases.

# Overfitting

Overfitting



$y$

$x_0$

- The learner fits the training data well, but loses the ability to generalize well, i.e. it has small training error but larger generalization error
- A learner with large capacity tends to overfit.
  - The family of functions is too large (compared with the size of the training data) and it contains many functions which all fit the training data well.
  - Without sufficient data, the learner cannot distinguish which one is most appropriate and would make an arbitrary choice among these apparently good solutions.
  - A separate validation set helps to choose a more appropriate one.
  - In most cases, data is contaminated by noise. The learner with large capacity tends to describe random errors or noise instead of the underlying models of data (classes).

# Model complexity (capacity)

- The goal is to classify novel examples not seen yet, but not the training examples!

- Generalization. The ability to correctly classify new examples that differ from those used for training



Overly complex models lead to complicated decision boundaries. It leads to perfect classification on the training examples, but would lead to poor performance on new patterns.

The decision boundary might represent the optimal tradeoff between performance on the training set and simplicity of classifier, therefore giving highest accuracy on new patterns.

# Generalization and Capacity



Typical relationship between capacity and both training and generalization (or test) error.

- As capacity increases, training error can be reduced, but the optimism (difference between training and generalization error) increases.

- At some point, the increase in optimism is larger than the decrease in training error (typically when the training error is low and cannot go much lower), and we enter the overfitting regime, where capacity is too large, above the optimal capacity.

- Before reaching optimal capacity, we are in the underfitting regime.

# Bias and Variance

Underfitting zone

Overfitting zone

Bias

Generalization error

Variance

Optimal capacity

Capacity

Image credit: Ian Goodfellow 160

# What about generalization ability of deep models?

- Deep nets have so many parameters they could just act like look up tables, regurgitating their training data

- Instead, they learn rules that generalize

- Defies classical theory!

A

under-fitting ⋮ over-fitting

Test risk

Risk

Training risk

sweet spot

Capacity of $\mathcal{H}$

B

under-parameterized ⋮ over-parameterized

Test risk

"classical" regime

"modern" interpolating regime

Risk

Training risk

interpolation threshold

Capacity of $\mathcal{H}$

[Double-descent: Belkin, Hsu, Ma, Mandal, PNAS 2019] 161

# Complexity and overfitting

- With limited training examples our polynomial regression model may achieve zero training error but nevertheless has a large test (generalization) error

train $$\frac{1}{n}\sum_{t=1}^{n}(y_t - f(x_t; \hat{\mathbf{w}}))^2 \approx 0$$

test $$E_{(x,y)\sim P}\left(y - f(x; \hat{\mathbf{w}})\right)^2 \gg 0$$



- We suffer from **overfitting** when the training error no longer bears any relation to the generalization error

# Avoiding overfitting: cross-validation

- **Cross-validation** allows us to estimate the generalization error based on training examples alone

Leave-one-out cross-validation treats each example in turn as a test example:

$$CV = \frac{1}{n} \sum_{i=1}^{n} \left( y_i - f(x_i; \hat{\mathbf{w}}^{-i}) \right)^2$$

where $\hat{\mathbf{w}}^{-i}$ are the least squares estimated parameters without the $i^{th}$ training example

# Polynomial regression: example (cont'd)



degree = 1, CV = 0.6    degree = 3, CV = 1.5

degree = 5, CV = 6.0    degree = 7, CV = 15.6

# Additive models

- More generally, predictions can be based on a linear combination of a set of basis functions (or features) $\{\phi_1(\mathbf{x}), \ldots, \phi_m(\mathbf{x})\}$, where each $\phi_i(\mathbf{x}) : \mathcal{R}^d \rightarrow \mathcal{R}$, and

$$f(\mathbf{x}; \mathbf{w}) = w_0 + w_1 \phi_1(\mathbf{x}) + \ldots + w_m \phi_m(\mathbf{x})$$

- Examples

If $\phi_i(x) = x^i$, $i = 1, \ldots, m$, then

$$f(x; \mathbf{w}) = w_0 + w_1 x + \ldots + w_{m-1} x^{m-1} + w_m x^m$$

# Additive models

- More generally, predictions can be based on a linear combination of a set of basis functions (or features) $\{\phi_1(\mathbf{x}), \dots, \phi_m(\mathbf{x})\}$, where each $\phi_i(\mathbf{x}) : \mathcal{R}^d \rightarrow \mathcal{R}$, and

$$f(\mathbf{x}; \mathbf{w}) = w_0 + w_1\phi_1(\mathbf{x}) + \dots + w_m\phi_m(\mathbf{x})$$

- Examples

If $\phi_i(x) = x^i$, $i = 1, \dots, m$, then

$$f(x; \mathbf{w}) = w_0 + w_1 x + \dots + w_{m-1} x^{m-1} + w_m x^m$$

If $m = d$, $\phi_i(\mathbf{x}) = x_i$, $i = 1, \dots, d$, then

$$f(\mathbf{x}; \mathbf{w}) = w_0 + w_1 x_1 + \dots + w_d x_d$$

# Additive models (cont'd)

- The basis functions can capture various (e.g., qualitative) properties of the inputs.

- For example: we can try to rate companies based on text descriptions

$$\mathbf{x} = \text{text document (collection of words)}$$

$$\phi_i(\mathbf{x}) = \begin{cases} 1 \text{ if word } i \text{ appears in the document} \\ 0 \text{ otherwise} \end{cases}$$

$$f(\mathbf{x}; \mathbf{w}) = w_0 + \sum_{i \in \mathsf{words}} w_i \phi_i(\mathbf{x})$$

# Additive models (cont'd)

- We can view the additive models graphically in terms of simple "units" and "weights"



- In **neural networks**, the basis functions themselves have adjustable parameters (cf. prototypes)

# Take-home messages



| Training error high? | → Yes → | Bigger model<br>Train longer<br>New model architecture | (Bias) |

**Training error high?** →Yes→ Bigger model / Train longer / New model architecture  (Bias)

↓ No

**Train-Dev error high?** →Yes→ More data / Regularization / New model architecture  (Variance)

↓ No

**Dev error high?** →Yes→ Make training data more similar to test data. / Data synthesis / (Domain adaptation.) / New model architecture  (Train-test data mismatch)

↓ No

**Test error high?** →Yes→ More dev set data  (Overfit dev set)

↓ No

Done!

# Statistical regression models

— model formulation, motivation
— maximum likelihood estimation

# Statistical view of linear regression

- In a statistical regression model we model both the function and noise

$$\textbf{Observed output} \quad = \quad \textbf{function} + \textbf{noise}$$

$$y \quad = \quad f(\mathbf{x}; \mathbf{w}) + \epsilon$$

where, e.g., $\epsilon \sim N(0, \sigma^2)$.

- Whatever we cannot capture with c
family of functions will be *interprete*

# Statistical view of linear regression

- $f(\mathbf{x}; \mathbf{w})$ is trying to capture the mean of the observations $y$ given the input $\mathbf{x}$:

$$E\{\, y \,|\, \mathbf{x}\,\} \;=\; E\{\, f(\mathbf{x}; \mathbf{w}) + \epsilon \,|\, \mathbf{x}\,\}$$
$$=\; f(\mathbf{x}; \mathbf{w})$$

where $E\{\, y \,|\, \mathbf{x}\,\}$ is the conditional expectation of $y$ given $\mathbf{x}$, evaluated according to the model (not according to the underlying distribution $P$)

# Statistical view of linear regression

- According to our statistical model

$$y = f(\mathbf{x}; \mathbf{w}) + \epsilon, \quad \epsilon \sim N(0, \sigma^2)$$

the outputs $y$ given $\mathbf{x}$ are normally distributed with mean $f(\mathbf{x}; \mathbf{w})$ and variance $\sigma^2$:

$$p(y|\mathbf{x}, \mathbf{w}, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\{ -\frac{1}{2\sigma^2}(y - f(\mathbf{x}; \mathbf{w}))^2 \}$$

(we model the uncertainty in the predictions, not just the mean)

- Loss function? Estimation?

# Maximum likelihood estimation

- Given observations $D_n = \{(\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_n, y_n)\}$ we find the parameters $\mathbf{w}$ that maximize the (conditional) likelihood of the outputs

$$L(D_n; \mathbf{w}, \sigma^2) = \prod_{i=1}^{n} p(y_i | \mathbf{x}_i, \mathbf{w}, \sigma^2)$$

- Example: linear function

$$p(y | \mathbf{x}, \mathbf{w}, \sigma^2) =$$

$$\frac{1}{\sqrt{2\pi\sigma^2}} \exp\{ -\frac{1}{2\sigma^2}(y - w_0 -$$

(why is this a bad fit according to the likelihood criterion?)

# Maximum likelihood estimation (cont'd)

Likelihood of the observed outputs:

$$L(D; \mathbf{w}, \sigma^2) = \prod_{i=1}^{n} P(y_i | \mathbf{x}_i, \mathbf{w}, \sigma^2)$$

- It is often easier (but equivalent) to try to maximize the log-likelihood:

$$
\begin{aligned}
l(D; \mathbf{w}, \sigma^2) &= \log L(D; \mathbf{w}, \sigma^2) = \sum_{i=1}^{n} \log P(y_i | \mathbf{x}_i, \mathbf{w}, \sigma^2) \\
&= \sum_{i=1}^{n} \left( -\frac{1}{2\sigma^2} (y_i - f(\mathbf{x}_i; \mathbf{w}))^2 - \log \sqrt{2\pi\sigma^2} \right) \\
&= \left( -\frac{1}{2\sigma^2} \right) \sum_{i=1}^{n} (y_i - f(\mathbf{x}_i; \mathbf{w}))^2 + \dots
\end{aligned}
$$

# Maximum likelihood estimation (cont'd)

- Maximizing log-likelihood is equivalent to minimizing empirical loss when the loss is defined according to

$$\text{Loss}(y_i, f(\mathbf{x}_i; \mathbf{w})) = -\log P(y_i | \mathbf{x}_i, \mathbf{w}, \sigma^2)$$

Loss defined as the negative log-probability is known as the log-loss.

# Maximum likelihood estimation (cont'd)

- The log-likelihood of observations

$$\log L(D; \mathbf{w}, \sigma^2) = \sum_{i=1}^{n} \log P(y_i | \mathbf{x}_i, \mathbf{w}, \sigma^2)$$

  is a generic fitting criterion and can be used to estimate the noise variance $\sigma^2$ as well.

- Let $\hat{\mathbf{w}}$ be the maximum likelihood (here least squares) setting of the parameters. What is the maximum likelihood estimate of $\sigma^2$, obtained by solving

$$\frac{\partial}{\partial \sigma^2} \log L(D; \mathbf{w}, \sigma^2) = 0 \quad ?$$

# Maximum likelihood estimation (cont'd)

- The log-likelihood of observations

$$\log L(D; \mathbf{w}, \sigma^2) = \sum_{i=1}^{n} \log P(y_i | \mathbf{x}_i, \mathbf{w}, \sigma^2)$$

  is a generic fitting criterion and can be used to estimate the noise variance $\sigma^2$ as well.

- Let $\hat{\mathbf{w}}$ be the maximum likelihood (here least squares) setting of the parameters. The maximum likelihood estimate of the noise variance $\sigma^2$ is

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^{n} (y_i - f(\mathbf{x}_i; \hat{\mathbf{w}}))^2$$

  i.e., the mean squared prediction error.

# Polynomial regression

- Consider again a simple $m^{th}$ degree polynomial regression model

$$y \;=\; w_0 + w_1 x + \ldots + w_m x^m + \epsilon, \;\; \epsilon \sim N(0, \sigma^2)$$

  where $\sigma^2$ is assumed fixed (known).

- In this model the outputs $\{y_1, \ldots, y_n\}$ corresponding to any inputs $\{x_1, \ldots, x_n\}$ are generated according to

$$\mathbf{y} \;=\; \mathbf{X}\mathbf{w} + \mathbf{e}, \quad \text{where}$$

$$\mathbf{y} \;=\; \begin{bmatrix} y_1 \\ \cdots \\ y_n \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} 1 & x_1 & \ldots & x_1^m \\ \cdots & \cdots & \cdots & \\ 1 & x_n & \ldots & x_n^m \end{bmatrix}, \quad \mathbf{e} = \begin{bmatrix} \epsilon_1 \\ \cdots \\ \epsilon_n \end{bmatrix}$$

  and $\epsilon_i \sim N(0, \sigma^2)$, $i = 1, \ldots, n$.

# ML estimator, uncertainty

- We are interested in studying how the choice of inputs $\{x_1,\ldots,x_n\}$ or, equivalently, $\mathbf{X}$, affects the accuracy of our regression model

- Our model for the outputs $\{y_1,\ldots,y_n\}$ given $\mathbf{X}$ is

$$\mathbf{y} = \mathbf{Xw} + \mathbf{e}, \quad \mathbf{e} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$$

- We assume also that the training outputs are actually generated by a model in this class with some fixed but unknown parameters $\mathbf{w}^*$ (same $\sigma^2$)

$$\mathbf{y} = \mathbf{Xw}^* + \mathbf{e}, \quad \mathbf{e} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$$

- We can now ask, for a given $\mathbf{X}$, how accurately we are able to recover the ''true'' parameters $\mathbf{w}^*$

# ML estimator, uncertainty

- The ML estimator $\hat{\mathbf{w}}$ viewed here as a function of the outputs $\mathbf{y}$ for a fixed $\mathbf{X}$, is given by

$$\hat{\mathbf{w}} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$$

- We need to understand how $\hat{\mathbf{w}}$ varies in relation to $\mathbf{w}^*$ when the outputs are generated according to

$$\mathbf{y} = \mathbf{X}\mathbf{w}^* + \mathbf{e}, \quad \mathbf{e} \sim N(\mathbf{0}, \sigma^2\mathbf{I})$$

- In the absence of noise $\mathbf{e}$, the ML estimator would recover $\mathbf{w}^*$ exactly (with only minor constraints on $\mathbf{X}$)

$$\begin{aligned} \hat{\mathbf{w}} &= (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T(\mathbf{X}\mathbf{w}^*) \\ &= (\mathbf{X}^T\mathbf{X})^{-1}(\mathbf{X}^T\mathbf{X})\mathbf{w}^* \\ &= \mathbf{w}^* \end{aligned}$$

# ML estimator, uncertainty

- In the presence of noise we can still use the fact that $\mathbf{y} = \mathbf{X}\mathbf{w}^* + \mathbf{e}$ to simplify the parameter estimates

$$
\begin{aligned}
\hat{\mathbf{w}} &= (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y} \\
&= (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T(\mathbf{X}\mathbf{w}^* + \mathbf{e}) \\
&= (\mathbf{X}^T\mathbf{X})^{-1}(\mathbf{X}^T\mathbf{X})\mathbf{w}^* + (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{e} \\
&= \mathbf{w}^* + (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{e}
\end{aligned}
$$

- So the ML estimate is the correct parameter vector plus an estimate based purely on noise

# Next Lecture:

# Multi-layer Perceptrons