

# CS4487 Machine Learning

---

## Lecture 1: Introduction

---

- Course Information
  - Content
    - Lecture: Learning algorithm & Mathematical Proof
    - Tutorial: Use learning algorithms on small examples
    - HW: four => derive and analyze machine learning algorithms with mathematical tools
  - Grade:
    - Coursework(40)
      - 10(tutorial)
      - 30(assignment)
    - midterm(30)
    - final(30)
  - Book:
    - Pattern Recognition & Machine Learning (C M B)
    - Convex Optimization (S B & L V)
    - The Matrix Cookbook (K B P)
    - Linear Algebra and Its Applications
    - Deep Learning
  - Syllabus:
    - Classification, regression, clustering, dimensionality reduction and deep learning
    - Model selection, model evaluation, regularization, design of experiments
  - Prerequisites:
    - Linear Algebra
    - Calculus
    - Probability and Statistics
    - Optimization
    - Python
  - Some problem
    - derivative of matrix with respect to another
    - maximum likelihood estimation(MLE) & maximum a posteriori estimation(MAP), and their relationship
    - difference between gradient descent and stochastic gradient descent
    - singular value decomposition(SVD) and principal component analysis (PCA) and their relationship
- Introduction
  - type
    - By I/O
      - Supervised: Learning to predict (with a set of input and output)
      - Unsupervised: Learning to organize and represent

- By Parameter
    - Non-parametric: storing the training data
    - Parametric: Using an algorithm to adapt the parameters in a mathematical or statistical model given training data
  - Sample => feature => training=>learned model
  - AI > ML (Data) > DL (multi-layer)
  - Math
    - Vector Space (column vector)
    - Matrix
      - Square, diagonal, identity
      - Transpose, trace, inverse, **determinant**
      - **eigenvalue, eigenvectors**, orthogonality
    - Probability Concepts
      - Probability Distribution
      - Random Variable (X): Maps each element(get 2 balls) of the sample space to a value x(2) in X(called the range of the random variable)
      - Marginalization
      - Conditioning
      - Bayes Rules
      - Expectations
      - Classical Distributions (Bernoulli, Multinomial, Gaussian)
- 

## Lecture 2: KNN and Naive Bayes

### Classification

- Definition:

#### Definition: The Classification Task

Given a feature vector  $\mathbf{x} \in \mathbb{R}^n$  that describes an object that belongs to one of  $|\mathcal{Y}|$  classes from the set  $\mathcal{Y}$ , predict which class the object belongs to

- Classifier Learning

#### Definition: Classifier Learning

Given a data set of example pairs  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)}), i = 1, \dots, m\}$  where  $\mathbf{x}^{(i)} \in \mathbb{R}^n$  is a feature vector and  $y^{(i)} \in \mathcal{Y}$  is a class label, learn a function  $f : \mathbb{R}^n \mapsto \mathcal{Y}$  that accurately predicts the class label  $y$  for any feature vector  $\mathbf{x}$

- Error and Accuracy

## Definition: Classification Error Rate

Given a data set of example pairs  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)}), i = 1, \dots, m\}$  and a function  $f : \mathbb{R}^n \mapsto \mathcal{Y}$ , the classification error rate of  $f$  on  $\mathcal{D}$  is

$$\text{Err}(f, \mathcal{D}) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}[y^{(i)} \neq f(\mathbf{x}^{(i)})]$$

## Definition: Classification Accuracy Rate

Given a data set of example pairs  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)}), i = 1, \dots, m\}$  and a function  $f : \mathbb{R}^n \mapsto \mathcal{Y}$ , the classification accuracy rate of  $f$  on  $\mathcal{D}$  is

$$\text{Acc}(f, \mathcal{D}) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}[y^{(i)} = f(\mathbf{x}^{(i)})]$$

## KNN

- Definition

The KNN classifier is a non-parametric classifier that simply stores the training data  $\mathcal{D}$  and classifies each new instance  $\mathbf{x}$  using a majority vote over its set of  $k$  nearest neighbors  $\mathcal{N}_k(\mathbf{x})$  computed using any distance function  $d : \mathbb{R}^n \times \mathbb{R}^n \mapsto \mathbb{R}$

## KNN Classification Function

$$f_{\text{KNN}}(\mathbf{x}) = \arg \max_{y \in \mathcal{Y}} \sum_{i \in \mathcal{N}_k(\mathbf{x})} \mathbb{I}[y^{(i)} = y]$$

Use of KNN requires choosing the distance function  $d$  and the number of neighbors  $k$

- Notice that *argmax* will traverse all the element in a set
- Work with any distance function  $d$  satisfying non-negativity  $d(x, x') > 0$  and identity of indiscernibles  $d(x, x) = 0$
- Alternatively, KNN can work with any similarity function  $s$  satisfying non-negativity  $s(x, x') \geq 0$  that attains its maximum on indiscernibles  $s(x, x) = \max_{x'} s(x, x')$
- Distance Metrics

## Definition: Minkowski Distance ( $\ell_p$ -norms)

Given two data vectors  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^n$ , the Minkowski Distance with parameter  $p$  (the  $\ell_p$ -norm) is a proper metric defined as follows:

$$\begin{aligned} d_p(\mathbf{x}, \mathbf{x}') &= \|\mathbf{x} - \mathbf{x}'\|_p \\ &= \left( \sum_{j=1}^n |x_j - x'_j|^p \right)^{1/p} \end{aligned}$$

Special cases include Euclidean distance ( $p = 2$ ), Manhattan distance ( $p = 1$ ) and Chebyshev distance ( $p = \infty$ )

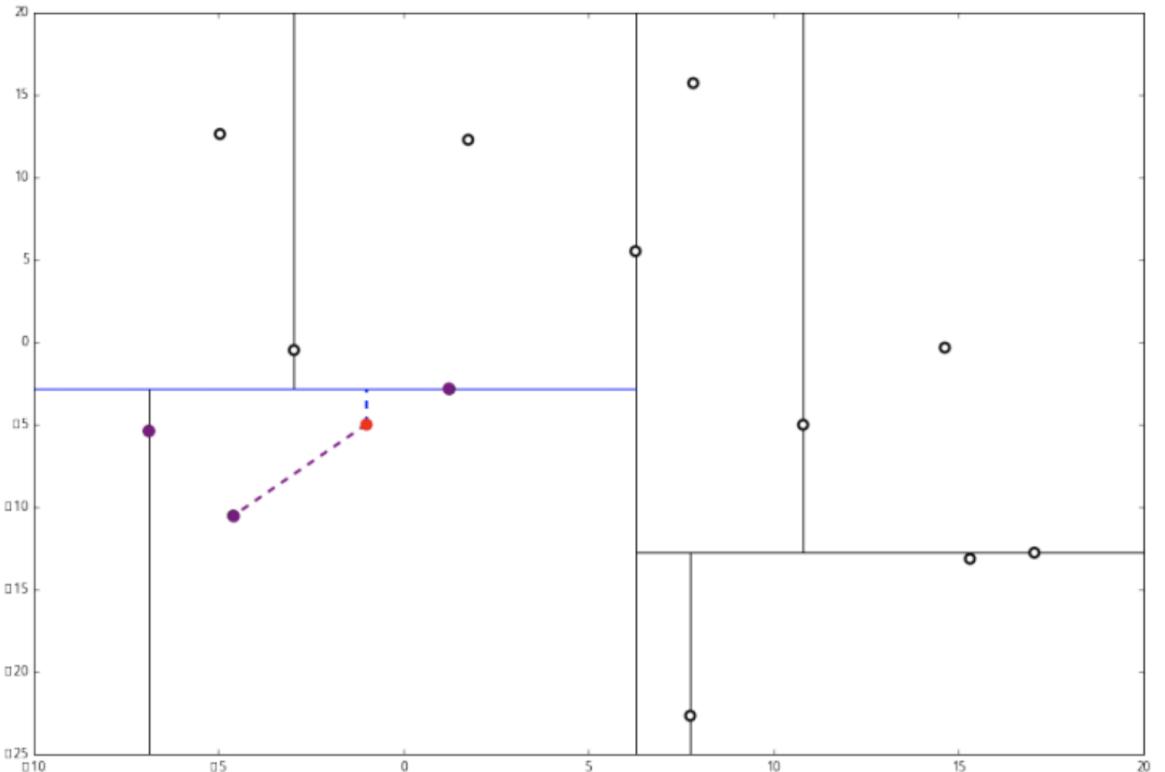
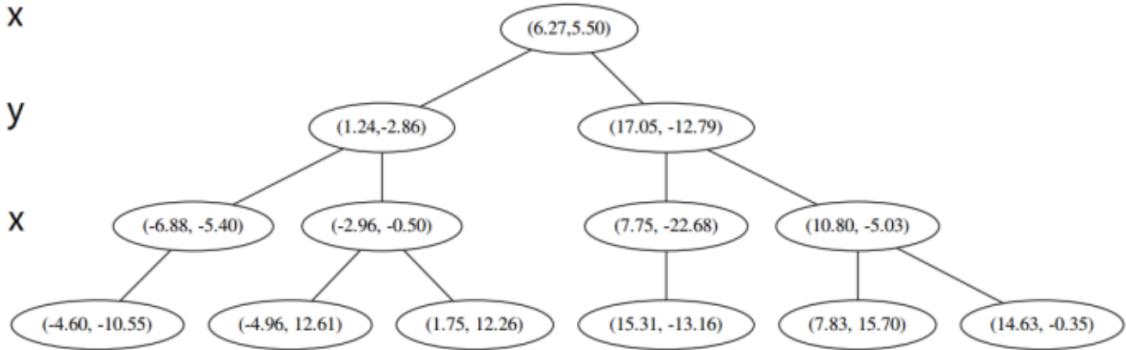
- Variant

- If there are 4 very far point and 1 very close point => vote to very far
- For each class  $y$ , traverse all the  $k$  closest points, add all weights that belongs to the class and divide by some of weight of all  $k$  closest points, where the weight of  $i_{th}$  point is defined as  $\exp(-\alpha d_i)$ . Then, choose the class with the maximum value
- Instead of giving all of the  $k$  neighbors equal weight in the majority vote, a distance-weighted majority can be used:

$$f_{\text{KNN}}(\mathbf{x}) = \arg \max_{y \in \mathcal{Y}} \frac{\sum_{i \in \mathcal{N}_k(\mathbf{x})} w_i \mathbb{I}[y^{(i)} = y]}{\sum_{i \in \mathcal{N}_k(\mathbf{x})} w_i},$$

$$w_i = \exp(-\alpha d_i)$$

- Instead of a brute force nearest neighbor search, data structures like  $k$ -d trees can be constructed over the training data that support nearest neighbor search with lower computational complexity
- [K-d tree to reduce the complexity](#): Make use of the division line
  - Point to search: (-1, -5)



- Trade off
  - Low bias: Converges to the correct decision surface as data goes to infinity
  - High variance: Lots of variability in the decision surface when amount of data is low
  - Curse of dimensionality: Everything is far from everything else in high dimensions
  - Space and time complexity: Need to store all training data. Take time at run time to compute

## Bayes Optimal Classifiers

- Probabilistic Classification
  - True probability of seeing a data case that belongs to  $p(Y = y) = \phi_y$
  - True probability density of seeing a data vector  $x \in \mathbb{R}^n$  that belongs to class  $y$  is  $p(X = x|Y = y) = \varphi_y(x)$
  - To calculate with vector  $x$  the probability of any one of  $y$  :

$$p(Y = y|X = x) = \frac{p(X = x|Y = y)p(Y = y)}{\sum_{\tilde{y} \in Y} p(X = x|Y = \tilde{y})p(Y = \tilde{y})} = \frac{\varphi_y(x)\phi_y}{\sum_{\tilde{y} \in Y} \varphi_{\tilde{y}}(x)\phi_{\tilde{y}}}$$

- Bayes Optimal Classifiers

$$f_B(x) = \operatorname{argmax}_{y \in Y} p(Y = y | X = x) = \operatorname{argmax}_{y \in Y} \varphi_y(x) \phi_y$$

$$\text{MIN Error} = 1 - E_x[\max_{y \in Y} p(Y = y | X = x)]$$

- $p(X = x)$  can be treated as a constant for specific input vector
- $E_x$  represents expectation with the condition  $X = x$
- Proof:

Derive the Bayes Optimal Classifier.

prof. ① Give classifier  $f$ . write down the 0-1 loss

for  $\underset{a}{\left[ \begin{array}{c} \text{training sample} \\ a \end{array} \right]} \quad \left[ \begin{array}{c} (x, y) \end{array} \right]$

$$\ell(f(x), y) = \mathbb{I}[f(x) \neq y] = 1 - \mathbb{I}[f(x) = y].$$

② Expected predicted Error:

$$\ell = \mathbb{E}[\ell(f(x), y)]$$

$$= \sum_{x \in X} \sum_{y \in Y} P(x=x, Y=y) \cdot \ell(f(x), y)$$

$$= \sum_{x \in X} P(x=x) \sum_{y \in Y} P(Y=y | x=x) \ell(f(x), y)$$

$$= \mathbb{E}_x \mathbb{E}_{y|x} [\ell(f(x), y)]$$

$$= \mathbb{E}_x [P(Y=y | x=x) \ell(f(x), y)]$$

③ The optimal classifier  $f_B$  should have the following condition for each  $X=x$

$$\begin{aligned}
 f_B(x) &= \underset{f}{\operatorname{argmax}} \sum_{y \in Y} \ell(f(x), y) P(Y=y | X=x) \\
 &= \underset{f}{\operatorname{argmin}} \sum_{y \in Y} [1 - \mathbb{I}[f(x)=y]] P(Y=y | X=x) \\
 &= \underset{f}{\operatorname{argmin}} \left( \sum_{y \in Y} P(Y=y | X=x) - \sum_{y \in Y} \mathbb{I}[f(x)=y] P(Y=y | X=x) \right) \\
 &= \underset{f}{\operatorname{argmin}} \left[ 1 - \sum_{y \in Y} \mathbb{I}[f(x)=y] P(Y=y | X=x) \right] \\
 &= \underset{\substack{\text{arg max } \\ y \in Y}}{\operatorname{arg max}} \sum_{y \in Y} \mathbb{I}[f(x)=y] P(Y=y | X=x) \quad \xrightarrow{\text{the classifier that } \sum_{y \in Y} \mathbb{I}[f(x)=y] P(Y=y | X=x) \text{ is fixed.}} \\
 &= \underset{\substack{\text{arg max } \\ y \in Y}}{\operatorname{arg max}} P(Y=y | X=x) \quad \xrightarrow{\text{选择最大的可能}}
 \end{aligned}$$

classifier chosen

Error for specific  $X=x$ :

$$\begin{aligned}
 &1 - \max_{y \in Y} P(Y=y | X=x) \\
 &= 1 - \mathbb{E}_{Y|X=x} [\max_{y \in Y} P(Y=y | X=x)].
 \end{aligned}$$

The classifier:

$$\begin{aligned}
 f_B(x) &= \underset{y \in Y}{\operatorname{argmax}} P(Y=y | X=x) \\
 &= \underset{y \in Y}{\operatorname{argmax}} \frac{P(X=x | Y=y) \cdot P(Y=y)}{P(X=x)} \quad \text{fixed} \\
 &\Rightarrow \underset{y \in Y}{\operatorname{argmax}} P(X=x | Y=y) P(Y=y) \\
 &= \underset{y \in Y}{\operatorname{argmax}} \varphi_y(x) \quad \varphi_y(x) \neq y \\
 &\quad \downarrow \\
 &\quad \text{difficult to estimate. (limited sample).}
 \end{aligned}$$

- Not useful in practice: difficult to estimate  $\varphi_y(x)$  in high dimensions

# Lecture03: Naive Bayes, Linear Discriminant Analysis and Logistic Regression

## Naive Bayes

- Naive bayes classifier approximates the Bayes optimal classifier using a simple form for the functions  $\varphi_y(x)$
- Assumes that all of the data dimensions are statistically independent given the value of the class variable

$$\varphi_y(x) = p(X = x|Y = y) = \prod_{j=1}^n p(X_j = x_j|Y = y) = \prod_{j=1}^n \varphi_{j|y}(x_j)$$

- The general form for the classification function:

$$f_{NB}(x) = \operatorname{argmax}_{y \in Y} \phi_y \prod_{j=1}^n \varphi_{j|y}(x_j)$$

- The functions  $\varphi_{j|y}(x_j) = p(X_j = x_j|Y = y)$  are called *marginal class conditional distributions*
  - For real valued  $x_j$ ,  $p(X_j = x_j|Y = y)$  is typically modeled as a Gaussian density  
 $\varphi_{j|y}(x_j) = N(x_j; \mu_{j|y}, \sigma_{j|y}^2)$
  - For binary valued  $x_j$ ,  $p(X_j = x_j|Y = y)$  is typically modeled as a Bernoulli distribution  
 $\varphi_{j|y}(x_j) = \theta_{j|y}^{x_j} (1 - \theta_{j|y})^{(1-x_j)}$
  - For general categorical values  $x_j$ ,  $p(X_j = x_j|Y = y)$  is typically modeled as a categorical distribution  
 $\varphi_{j|y}(x_j) = \prod_{x \in X_j} \theta_{x,j|y}^{I[x_j=x]}$
- The problem to maximize is how to choose a correct  $\phi_y$  and  $\varphi_{j|y}(x_j) \Rightarrow$  using the *maximum likelihood* over  $D = \{(x^{(i)}, y^{(i)}), i = 1, \dots, m\}$ 
  - **posterior probability:**  $P(Y = y|X = x), \dots$  : likelihood  $P(X = x|Y = y)$ , Prior probability:  $P(Y = y)$
  - maximum likelihood: 让取样获得的大前提( $P(D)$ )可能性最大,  $D$  为数据集

- Class probabilities:  $\phi_y = \frac{1}{m} \sum_{i=1}^m \mathbb{I}[y^{(i)} = y]$
- Gaussian:  $\mu_{j|y} = \frac{\sum_{i=1}^m \mathbb{I}[y^{(i)}=y] x_j^{(i)}}{\sum_{i=1}^m \mathbb{I}[y^{(i)}=y]}, \quad \sigma_{j|y}^2 = \frac{\sum_{i=1}^m \mathbb{I}[y^{(i)}=y] (x_j^{(i)} - \mu_{j|y})^2}{\sum_{i=1}^m \mathbb{I}[y^{(i)}=y]}$
- Bernoulli:  $\theta_{j|y} = \frac{\sum_{i=1}^m \mathbb{I}[y^{(i)}=y] x_j^{(i)}}{\sum_{i=1}^m \mathbb{I}[y^{(i)}=y]}$
- Categorical:  $\theta_{x,j|y} = \frac{\sum_{i=1}^m \mathbb{I}[y^{(i)}=y \cap x_j^{(i)}=x]}{\sum_{i=1}^m \mathbb{I}[y^{(i)}=y]}$

- for Bernoulli distribution,  $x_j$  is 0 or 1
- proof of the class probabilities  $\phi_y$  and the categorical probability  $\theta_{x,j|y} \Rightarrow$  see lecture note
- Geometric Interpretation: for normal distribution of  $x$ , and binary classification problem we have  $\Rightarrow$

$$\begin{aligned}
f_{\text{NB}}(\mathbf{x}) &= \arg \max_{y \in \mathcal{Y}} \phi_y \prod_{j=1}^n \varphi_{j|y}(x_j) \\
&= \arg \max_{y \in \mathcal{Y}} \log(\phi_y) + \sum_{j=1}^n \log(\varphi_{j|y}(x_j)) \\
&= \arg \max_{y \in \mathcal{Y}} \log(\phi_y) + \sum_{j=1}^n \left( -\frac{1}{2} \log(2\pi\sigma_{j|y}^2) - \frac{(x_j - \mu_{j|y})^2}{2\sigma_{j|y}^2} \right)
\end{aligned}$$

- The decision boundary consists of the set of points  $\mathbf{x}$  where:

$$\begin{aligned}
&\log(\phi_0) + \sum_{j=1}^n \left( -\frac{1}{2} \log(2\pi\sigma_{j|0}^2) - \frac{1}{2\sigma_{j|0}^2} (x_j - \mu_{j|0})^2 \right) \\
&- \log(\phi_1) - \sum_{j=1}^n \left( -\frac{1}{2} \log(2\pi\sigma_{j|1}^2) - \frac{1}{2\sigma_{j|1}^2} (x_j - \mu_{j|1})^2 \right) = 0
\end{aligned}$$

- It's easy to see that the decision boundary is a quadratic function of  $\mathbf{x}$  with the form:  $\sum_{j=1}^n (a_j x_j^2 + b_j x_j) + c = 0$
- In the multi-class case, the decision boundary is piece-wise quadratic

- Laplace Smoothing
  - The maximum likelihood :

$$\therefore \theta_{x,j|y} = \frac{\sum_{i=1}^m \mathbb{I}[y^{(i)}=y \cap x_j^{(i)}=x]}{\sum_{i=1}^m \mathbb{I}[y^{(i)}=y]}$$

- If the training set is small, it's possible that  $I[x_j^{(i)} = x] = 0$ , for any  $i \in [1, m]$ ,  $i$  is integer
- Cause the probability to be zero, which is overfitting
- Laplace smoothing method:

$$\theta_{x,j|y} = \frac{\sum_{i=1}^m \mathbb{I}[y^{(i)} = y \cap x_j^{(i)} = x] + \alpha}{\sum_{i=1}^m \mathbb{I}[y^{(i)} = y] + \alpha |\mathcal{X}|},$$

- When  $\alpha = 1$  it's called Laplace smoothing
- it can be shown that Laplace smoothing is the posterior mean with Dirichlet distribution as the conjugate prior for Categorical distribution
- Trade offs
  - Speed: learning and classification are fast
  - Storage:  $O(n)$  parameters, large compression of the training data
  - Interpretability: good because  $\varphi_y(x)$  is the class conditional averages
  - Accuracy: feature independence assumption and canonical forms for class-conditional marginals will rarely be correct for real-world problems, leading to lower accuracy
  - Data: Some care is needed in the estimation of parameters in the discrete case when data is scarce(small data set)

## Linear Discriminant Analysis

- Linear Discriminant Analysis
  - A different approximation to the Bayes optimal classifier for real-valued data (instead of Maximum Likelihood)
  - Use dependent Normal Distribution in Naive Bayes, LDA assumes  $\varphi_y(x) = N(x; \mu_y, \Sigma)$ , where  $\Sigma$  is covariance matrix

$$\varphi_y(\mathbf{x}) = \frac{1}{|(2\pi)^n \Sigma|^{1/2}} \exp \left( -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_y)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_y) \right)$$

- The classification function is  $f_{\text{LDA}}(\mathbf{x}) = \arg \max_{y \in \mathcal{Y}} \varphi_y(\mathbf{x}) \phi_y$

To get an intuition for what a multivariate Gaussian is, consider the simple case where  $n = 2$ , and where the covariance matrix  $\Sigma$  is diagonal, i.e.,

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad \mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} \quad \Sigma = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}$$

In this case, the multivariate Gaussian density has the form,

$$\begin{aligned} p(x; \mu, \Sigma) &= \frac{1}{2\pi \left| \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix} \right|^{1/2}} \exp \left( -\frac{1}{2} \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{bmatrix}^T \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}^{-1} \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{bmatrix} \right) \\ &= \frac{1}{2\pi(\sigma_1^2 \cdot \sigma_2^2 - 0 \cdot 0)^{1/2}} \exp \left( -\frac{1}{2} \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{bmatrix}^T \begin{bmatrix} \frac{1}{\sigma_1^2} & 0 \\ 0 & \frac{1}{\sigma_2^2} \end{bmatrix} \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{bmatrix} \right), \end{aligned}$$

where we have relied on the explicit formula for the determinant of a  $2 \times 2$  matrix<sup>3</sup>, and the fact that the inverse of a diagonal matrix is simply found by taking the reciprocal of each diagonal entry. Continuing,

$$\begin{aligned} p(x; \mu, \Sigma) &= \frac{1}{2\pi\sigma_1\sigma_2} \exp \left( -\frac{1}{2} \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{bmatrix}^T \begin{bmatrix} \frac{1}{\sigma_1^2}(x_1 - \mu_1) \\ \frac{1}{\sigma_2^2}(x_2 - \mu_2) \end{bmatrix} \right) \\ &= \frac{1}{2\pi\sigma_1\sigma_2} \exp \left( -\frac{1}{2\sigma_1^2}(x_1 - \mu_1)^2 - \frac{1}{2\sigma_2^2}(x_2 - \mu_2)^2 \right) \\ &= \frac{1}{\sqrt{2\pi}\sigma_1} \exp \left( -\frac{1}{2\sigma_1^2}(x_1 - \mu_1)^2 \right) \cdot \frac{1}{\sqrt{2\pi}\sigma_2} \exp \left( -\frac{1}{2\sigma_2^2}(x_2 - \mu_2)^2 \right). \end{aligned}$$

- By using maximum likelihood, which reduces to using sample estimates:

- Class probabilities:  $\phi_y = \frac{1}{m} \sum_{i=1}^m \mathbb{I}[y^{(i)} = y]$
- Class means:  $\boldsymbol{\mu}_y = \frac{\sum_{i=1}^m \mathbb{I}[y^{(i)} = y] \mathbf{x}^{(i)}}{\sum_{i=1}^m \mathbb{I}[y^{(i)} = y]}$
- Shared covariance:  $\boldsymbol{\Sigma} = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}^{(i)} - \boldsymbol{\mu}_{y^{(i)}})(\mathbf{x}^{(i)} - \boldsymbol{\mu}_{y^{(i)}})^T$

- similar to the one dimensional case, the decision boundary consists of points  $x$  where:

$$\begin{aligned} \log(\phi_0) - \frac{1}{2} \log |(2\pi)^n \boldsymbol{\Sigma}| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_0) \\ - \log(\phi_1) + \frac{1}{2} \log |(2\pi)^n \boldsymbol{\Sigma}| + \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) = 0 \end{aligned}$$

- We can cancel a large number of terms because of the common covariance matrix and obtain the following result:

$$\log(\phi_0) - \log(\phi_1) - 0.5 \boldsymbol{\mu}_0^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_0 + 0.5 \boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 + (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} \mathbf{x} = 0$$

- Trade-offs

- Speed: quadratic dependence on  $n$  dimension makes LDA slower than Naive Bayes
- Storage:  $O(n^2)$  parameters, good when  $n \ll m$
- Interpretability: Good interpretability
- Accuracy: Rarely be correct in real-world problems. However, the induced linear decision boundaries can often perform reasonably well
- Data: LDA will generally need more data than NB since it needs to estimate the  $O(n^2)$  parameters in the pooled covariance matrix

## Logistic Regression

- Generative vs. Discriminative Classifiers
  - Naive Bayes and LDA are said to be *generative* classifier because they explicitly model the joint distribution  $p(x, y)$  of the data vectors  $x$  and the label  $y$
  - To build a probabilistic classifier, all we really need to model is  $p(y|x)$
  - classifiers based on directly estimating  $p(y|x)$  are called *discriminative* classifiers because they ignore the distribution of  $x$  and focus only on the class labels  $y$
- Logistic Regression
  - In binary case, it directly models the decision boundary using a linear function:

$$\log p(Y = 1|\mathbf{x}) - \log p(Y = 0|\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$$

$$p(Y = 1|\mathbf{x}) = \frac{1}{1 + \exp(-(\mathbf{w}^T \mathbf{x} + b))}$$

- Classification function based on  $\mathbf{x}$  (not  $\mathbf{y}$ )*
- Now the classification function is :  $f_{LR} = \text{argmax}_{y \in Y} p(Y = y|x)$ :
  - the standard logistic function, or the sigmoid function is defined as:

$$\begin{aligned} g(z) &= \frac{1}{1 + e^{-z}} \\ g'(z) &= -\left(\frac{1}{1 + e^{-z}}\right)^2 * (-e^{-z}) \\ &= g(z)(1 - g(z)) \end{aligned}$$

- The learning logistic regression
  - logistic regression model parameters  $\theta = \{w, b\}$  are selected to optimize the conditional log likelihood of the labels given the data set  $D = \{(x^{(i)}, y^{(i)}), i = 1, \dots, m\}$ :

$$\begin{aligned} \boldsymbol{\theta}^* &= \arg \max_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta} | D) = \arg \max_{\boldsymbol{\theta}} \frac{1}{m} \sum_{i=1}^m \log p(Y = y^{(i)} | \mathbf{X} = \mathbf{x}^{(i)}) \\ &= \arg \max_{\boldsymbol{\theta}} \frac{1}{m} \sum_{i=1}^m y^{(i)} \log g(\boldsymbol{\theta}^T \mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - g(\boldsymbol{\theta}^T \mathbf{x}^{(i)})) \end{aligned}$$

- cannot be maximized analytically. Learning the model parameters requires numerical optimization methods
- Multiclass Logistic Regression
  - Logistic regression extends to multiclass case:

$$p(Y = y|\mathbf{x}) = \frac{\exp(\mathbf{W}_y \mathbf{x} + b_y)}{\sum_{\tilde{y} \in \mathcal{Y}} \exp(\mathbf{W}_{\tilde{y}} \mathbf{x} + b_{\tilde{y}})}$$

- [why different from binary case?](#)
- Geometry
  - **Explicitly designed** to have a **linear decision boundary in the binary case**
  - In the multiclass case, the decision boundary is piece-wise linear
  - Same representational capacity as LDA
- Trade-offs
  - Speed: Faster than NB and LDA in classification time, numerical optimization is slower than naive bayes
  - Storage:  $O(n)$  parameters
  - Interpretability: The “importance” of different feature variables  $x^j$  can be understood in terms of their weights  $w^j$
  - Accuracy: Tends to be better in high dimensions with limited data compared to LDA. Much worse than KNN in low dimensions with **lots of data and non-linear decision boundaries**
- Probability view of logistic function(sigmoid function)

## 2.2 Probabilistic View on the Sigmoid Function

Let's write the general form of the sigmoid function

$$F(z; \mu, s) = \frac{1}{1 + e^{-\frac{z-\mu}{s}}}, \quad (25)$$

where  $\mu$  and  $s > 0$  are location and scale parameters, respectively.

For a function  $F(z)$  to be a legitimate cumulative distribution function (CDF)<sup>4</sup>, it must satisfy four properties:

- Non-decreasing;
- Right-continuous (*i.e.*,  $\lim_{z \rightarrow c^+} F(z) = F(c)$ );
- $\lim_{z \rightarrow -\infty} F(z) = 0$ ;
- $\lim_{z \rightarrow +\infty} F(z) = 1$ .

It is not hard to see that the sigmoid function satisfies all of the properties, and the corresponding probability density function (PDF) is

$$f(z) = F'(z) = \frac{e^{-r}}{s(1 + e^{-r})^2} \quad \text{and} \quad r = \frac{z - \mu}{s}, \quad (26)$$

which is known as the logistic distribution<sup>5</sup>.

# Lecture 04: Gradient Descent and Stochastic Gradient Descent

## Gradient Descent

- Gradient Descent
  - typical form of machine learning problem

$$\min_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \mathcal{D}) = \min_{\boldsymbol{\theta}} \frac{1}{m} \sum_{i=1}^m \ell(\boldsymbol{\theta}; \mathbf{x}^{(i)}, y^{(i)}),$$

where  $\ell : \mathbb{R}^n \mapsto \mathbb{R}$  is differentiable but not necessarily convex, and  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)}), i = 1, \dots, m\}$  is a finite training set

- for logistic regression, trying to minimize the negative log likelihood:

$$\min_{\boldsymbol{\theta}} -\frac{1}{m} \sum_{i=1}^m \log p(y^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta})$$

- choose initial point  $\theta^{(0)} \in \mathbb{R}^n$  and repeats

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \nabla \ell(\theta^{(t)}; \mathcal{D}) = \theta^{(t)} - \alpha \frac{1}{m} \sum_{i=1}^m \nabla \ell(\theta^{(t)}; \mathbf{x}^{(i)}, y^{(i)})$$

How to stop?  $\|\theta^{(t+1)} - \theta^{(t)}\| \leq \epsilon$  or  $\|\nabla \ell(\theta^{(t)}; \mathcal{D})\| \leq \epsilon$ , where  $\epsilon$  is a small positive constant, e.g.,  $10^{-6}$

- why gradient descent works

*Proof:* Directional derivative is along vector  $\vec{v} = [v_1, v_2, \dots, v_n]^T$  is defined by the function

$$\nabla_{\vec{v}} f(\vec{z}) = \lim_{\alpha \rightarrow 0} \frac{f(\vec{z} + \alpha \vec{v}) - f(\vec{z})}{\alpha}$$

Method 1:

$$\begin{aligned} \text{Suppose } \vec{x} &= \vec{z} + \alpha \vec{v} \\ \vec{x} &= \begin{bmatrix} z_1 + \alpha v_1 \\ z_2 + \alpha v_2 \\ \vdots \\ z_n + \alpha v_n \end{bmatrix} \end{aligned}$$

$$\begin{aligned} \frac{d f(\vec{z} + \alpha \vec{v})}{d \alpha} &= \sum_{i=1}^n \frac{df}{dx_i} \cdot \frac{dx_i}{d \alpha} \\ &= \sum_{i=1}^n \frac{\partial f}{\partial z_i} v_i \\ &= \nabla f(\vec{z})^T \vec{v} \end{aligned}$$

Method 2:

$$\begin{aligned} f(\vec{x}) &= \sum_{i=0}^n \frac{\nabla f^{(i)}(\vec{z})^T (\vec{x} - \vec{z})^{(i)}}{(i)!} + R(n) \\ &= f(\vec{z}) + \nabla f'(\vec{z})^T (\vec{x} - \vec{z}) + \sum_{i=2}^n \frac{\nabla f^{(i)}(\vec{z})^T (\vec{x} - \vec{z})^{(i)}}{(i)!} + R(n) \\ \hookrightarrow f(\vec{z} + \alpha \vec{v}) &= f(\vec{z}) + \alpha \nabla f'(\vec{z})^T \vec{v} + \sum_{i=2}^n \frac{\nabla f^{(i)}(\vec{z})^T (\vec{v})^{(i)}}{(i)!} \cdot \vec{v}^{(i)} + R(n) \\ &= f(\vec{z}) + \alpha \nabla f'(\vec{z})^T \vec{v} + O(\alpha^2), \\ &= f(\vec{z}) + \alpha \nabla f(\vec{z})^T \vec{v} + O(\alpha^2) - f(\vec{z}) \end{aligned}$$

$$\begin{aligned} \nabla_{\vec{v}} f &= \lim_{\alpha \rightarrow 0} \frac{f(\vec{z})^T \vec{v} + O(\alpha)}{\alpha} \\ &= \lim_{\alpha \rightarrow 0} \left( f(\vec{z})^T \vec{v} + O(\alpha) \right) \\ &= \nabla f(\vec{z})^T \vec{v} \end{aligned}$$

**Reason.** Pick any direction  $\mathbf{d}$  with  $\|\mathbf{d}\| = 1$ . The rate of change at  $\boldsymbol{\theta}$  is

$$\nabla \ell(\boldsymbol{\theta})^T \mathbf{d} \leq \|\nabla \ell(\boldsymbol{\theta})\| \|\mathbf{d}\| = \|\nabla \ell(\boldsymbol{\theta})\|.$$

If we set  $\mathbf{d} = \frac{\nabla \ell(\boldsymbol{\theta})}{\|\nabla \ell(\boldsymbol{\theta})\|}$ , then

$$\nabla \ell(\boldsymbol{\theta})^T \mathbf{d} = \|\nabla \ell(\boldsymbol{\theta})\|$$

- Moreover, a negative gradient step can decrease the loss function

**Reason.** Let  $\boldsymbol{\theta}^{(0)}$  be any initial point. Using the first order Taylor expansion for candidate point  $\boldsymbol{\theta}^{(1)} = \boldsymbol{\theta}^{(0)} - \alpha \nabla \ell(\boldsymbol{\theta}^{(0)})$ , we have

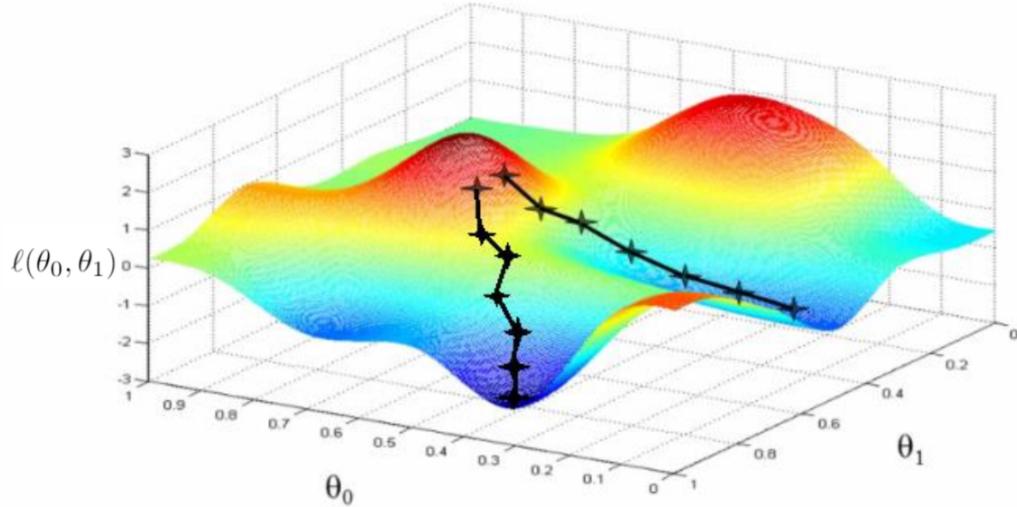
$$\ell(\boldsymbol{\theta}^{(1)}) - \ell(\boldsymbol{\theta}^{(0)}) = -\alpha \|\nabla \ell(\boldsymbol{\theta}^{(0)})\|^2 + O(\alpha^2).$$

Hence, if  $\nabla \ell(\boldsymbol{\theta}^{(0)}) \neq 0$  and  $\alpha$  is sufficiently small, we have

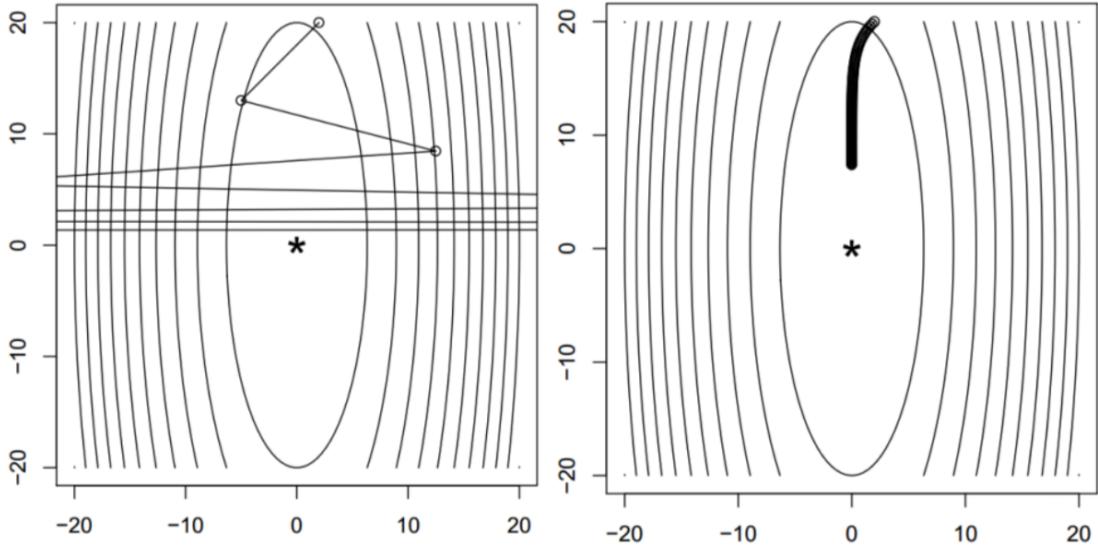
$$\ell(\boldsymbol{\theta}^{(1)}) < \ell(\boldsymbol{\theta}^{(0)}).$$

Therefore, for sufficiently small  $\alpha$ ,  $\boldsymbol{\theta}^{(1)}$  is an improvement over  $\boldsymbol{\theta}^{(0)}$

Initial point matters.



Learning rate matters.



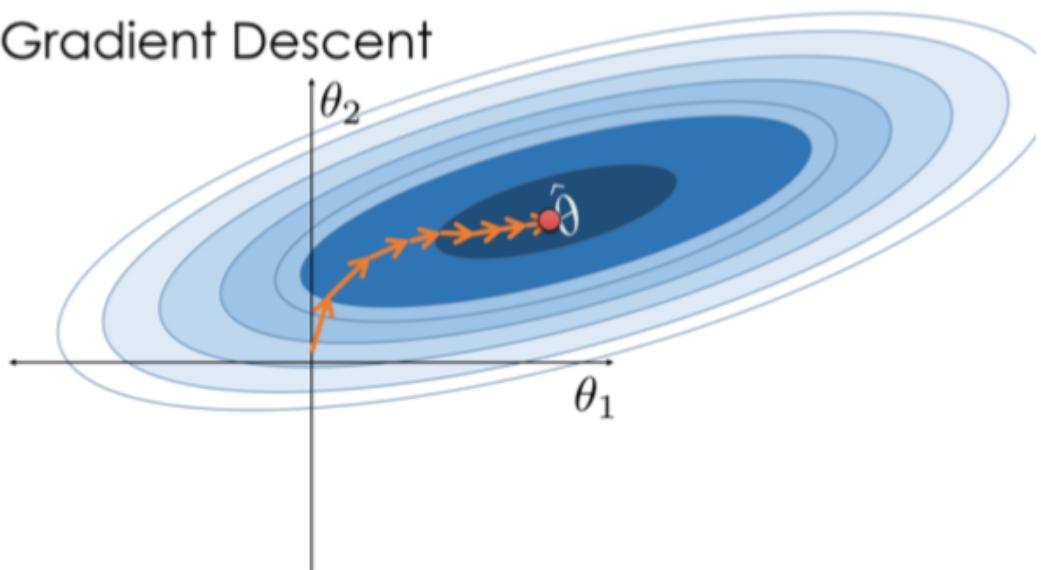
- Stochastic Gradient Descent
  - Compute loss for all the training data to get average takes  $O(m)$  times
  - Can only choose one training data(i.e. one vector) to (point)estimate the true gradient

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \alpha \nabla \ell(\boldsymbol{\theta}^{(t)}; \mathbf{x}^{(i)}, y^{(i)})$$

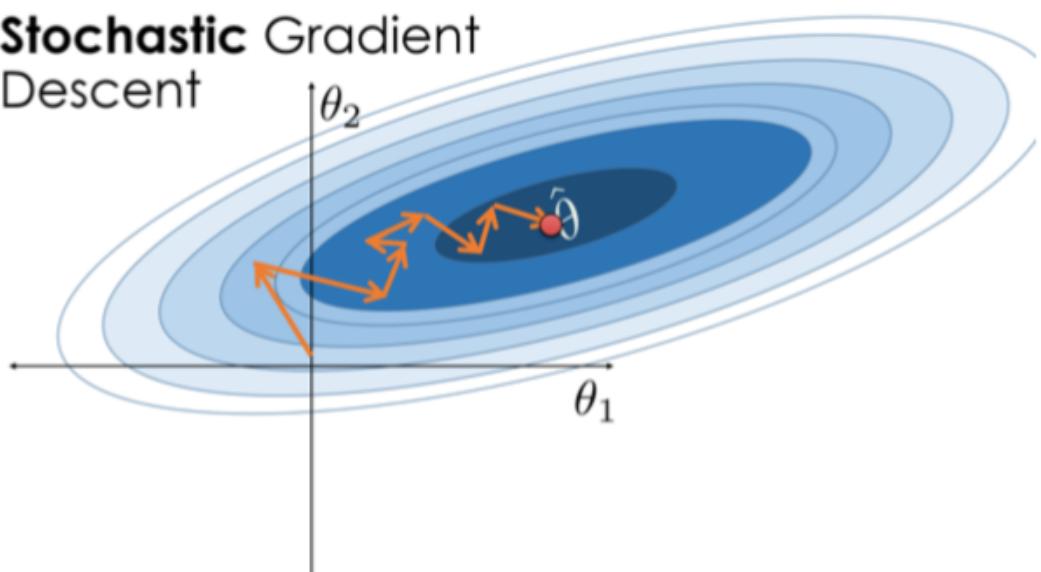
- The stochastic gradient  $\nabla \ell(\boldsymbol{\theta}; \mathbf{x}^{(i)}, y^{(i)})$  is the unbiased estimate of gradient  $\nabla \ell(\boldsymbol{\theta}; \mathcal{D})$

$$\mathbb{E}_i[\nabla \ell(\boldsymbol{\theta}; \mathbf{x}^{(i)}, y^{(i)})] = \frac{1}{m} \sum_{i=1}^m \nabla \ell(\boldsymbol{\theta}; \mathbf{x}^{(i)}, y^{(i)}) = \nabla \ell(\boldsymbol{\theta}; \mathcal{D})$$

Gradient Descent



**Stochastic** Gradient  
Descent



## Pros

- It is easy to fit into memory and computationally fast due to a single training sample being processed at each iteration
- The steps have oscillations, which may help the algorithm get out of bad local minima

## Cons

- The steps are noisy. The objective does not always decrease for each step and it may take longer to achieve convergence
- It loses the advantage of vectorized operations as it deals with only a single example at a time

- Mini-Batch Gradient Descent

- use a subset of samples => interval estimate => mean value theorem

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \alpha \nabla \ell(\boldsymbol{\theta}^{(t)}; \mathcal{B}) = \boldsymbol{\theta}^{(t)} - \alpha \frac{1}{|\mathcal{B}|} \sum_{(\mathbf{x}, y) \in \mathcal{B}} \nabla \ell(\boldsymbol{\theta}^{(t)}; \mathbf{x}, y)$$

- Mini-batch gradient descent and its variants are extremely popular and extensively used in machine learning

- Gradient Descent of Logistic Regression

### 3.2 Gradient Descent for Logistic Regression

*Proof.* Let's start by working with just one training example  $(x, y)$ , and take derivatives to derive the stochastic gradient ascent rule:

$$\begin{aligned}
 \frac{\partial}{\partial \theta_j} \ell(\theta) &= \left( y \frac{1}{g(\theta^T x)} - (1-y) \frac{1}{1-g(\theta^T x)} \right) \frac{\partial}{\partial \theta_j} g(\theta^T x) \\
 &= \left( y \frac{1}{g(\theta^T x)} - (1-y) \frac{1}{1-g(\theta^T x)} \right) g(\theta^T x)(1-g(\theta^T x)) \frac{\partial}{\partial \theta_j} \theta^T x \\
 &= (y(1-g(\theta^T x)) - (1-y)g(\theta^T x)) x_j \\
 &= (y - g(\theta^T x)) x_j.
 \end{aligned} \tag{30}$$

This therefore gives us the stochastic gradient ascent rule

$$\theta_j^{(t+1)} = \theta_j^{(t)} + \alpha(y^{(i)} - g((\theta^{(t)})^T x^{(i)}))x_j^{(i)}, \tag{31}$$

from which we can easily derive the gradient ascent rule

$$\theta_j^{(t+1)} = \theta_j^{(t)} + \alpha \frac{1}{m} \sum_{i=1}^m (y^{(i)} - g((\theta^{(t)})^T x^{(i)}))x_j^{(i)}. \tag{32}$$

Finally, the corresponding vectorized gradient ascent rule is

$$\theta^{(t+1)} = \theta^{(t)} + \alpha \frac{1}{m} \sum_{i=1}^m (y^{(i)} - g((\theta^{(t)})^T x^{(i)}))x^{(i)}. \tag{33}$$

□

## Lecture05: Linear Regression

### Regression

- Regression: Given a feature vector  $x \in \mathbb{R}^n$ , predict its corresponding output value  $y \in \mathbb{R}$
- Questions:
  - Does size of sample matters? => we prefer larger sample size
  - Does location of points matters? => Samples need be representative or informative
  - Given a set of  $m$  points, is curve fitting unique? => Estimating “continuous” model from “discrete” data is all ill-posed problem and can never be “certain”
  - How to choose the best model from infinite curves? => Have some criteria choose preferred model
- Example: House price, weather prediction, age estimation ...
- Regression learning problem:

#### Definition: Regression Learning Problem

Given a data set of example pairs  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)}), i = 1, \dots, m\}$  where  $\mathbf{x}^{(i)} \in \mathbb{R}^n$  is a feature vector and  $y^{(i)} \in \mathbb{R}$  is the output, learn a function  $f : \mathbb{R}^n \mapsto \mathbb{R}$  that accurately predicts  $y$  for any feature vector  $\mathbf{x}$

- Error measure:

- MSE:

### Definition: Mean Squared Error (MSE)

Given a data set of example pairs  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)}), i = 1, \dots, m\}$  and a function  $f : \mathbb{R}^n \mapsto \mathbb{R}$ , the MSE of  $f$  on  $\mathcal{D}$  is:

$$\text{MSE}(f, \mathcal{D}) = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - f(\mathbf{x}^{(i)}))^2$$

Related measures include:

Sum of Squared Errors:  $\text{SSE}(f, \mathcal{D}) = m \cdot \text{MSE}(f, \mathcal{D})$

Risidual Sum of Squares:  $\text{RSS}(f, \mathcal{D}) = m \cdot \text{MSE}(f, \mathcal{D})$

Root Mean Squared Error:  $\text{RMSE}(f, \mathcal{D}) = \sqrt{\text{MSE}(f, \mathcal{D})}$

- SSE and RSS are the same

- MAE:

### Definition: Mean Absolute Error (MAE)

Given a data set of example pairs  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)}), i = 1, \dots, m\}$  and a function  $f : \mathbb{R}^n \mapsto \mathbb{R}$ , the MAE of  $f$  on  $\mathcal{D}$  is:

$$\text{MAE}(f, \mathcal{D}) = \frac{1}{m} \sum_{i=1}^m |y^{(i)} - f(\mathbf{x}^{(i)})|$$

## Linear Regression

- A parametric regression method that assumes the relationship between  $y$  and  $x$  is a linear function with parameters  $w = [w_1, w_2, \dots, w_n]^T$  and  $b$

$$f_{\text{Lin}}(\mathbf{x}) = \left( \sum_{j=1}^n w_j x_j \right) + b = \mathbf{x}^T \mathbf{w} + b$$

- Ordinary Least Squares (OLS) Linear Regression

- Minimize the MSE on the training data set:

$$\mathbf{w}^*, b^* = \arg \min_{\mathbf{w}, b} \frac{1}{m} \sum_{i=1}^m (y^{(i)} - (\mathbf{x}^{(i)})^T \mathbf{w} - b)^2$$

- Use derivative:

for two-dimension case:

$$\begin{aligned}
 & \underset{\mathbf{w}, b}{\operatorname{arg} \min} \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2 \\
 &= \underset{\mathbf{w}, b}{\operatorname{arg} \min} \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \mathbf{w} \mathbf{x}^{(i)} - b)^2 \\
 & J(\mathbf{w}, b) = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \mathbf{w} \mathbf{x}^{(i)} - b)^2 \\
 & \frac{\partial J}{\partial \mathbf{w}} = -\frac{2}{m} \sum_{i=1}^m (y^{(i)} - \mathbf{w} \mathbf{x}^{(i)} - b) \mathbf{x}^{(i)} = 0 \\
 & \frac{\partial J}{\partial b} = -\frac{2}{m} \sum_{i=1}^m (y^{(i)} - \mathbf{w} \mathbf{x}^{(i)} - b) = 0 \\
 & \left\{ \begin{array}{l} \mathbf{w} \sum_{i=1}^m (\mathbf{x}^{(i)})^T + b \sum_{i=1}^m \mathbf{x}^{(i)} = \sum_{i=1}^m (\mathbf{x}^{(i)} y^{(i)}) \\ \mathbf{w} \sum_{i=1}^m \mathbf{x}^{(i)} + b m = \sum_{i=1}^m y^{(i)} \end{array} \right. \\
 & \left[ \begin{array}{cc} \sum_{i=1}^m (\mathbf{x}^{(i)})^T & \sum_{i=1}^m \mathbf{x}^{(i)} \\ \sum_{i=1}^m \mathbf{x}^{(i)} & m \end{array} \right] \left[ \begin{array}{c} \mathbf{w} \\ b \end{array} \right] = \left[ \begin{array}{c} \sum_{i=1}^m \mathbf{x}^{(i)} y^{(i)} \\ \sum_{i=1}^m y^{(i)} \end{array} \right] \\
 & \left[ \begin{array}{c} \mathbf{w} \\ b \end{array} \right] = \left[ \begin{array}{cc} \sum_{i=1}^m (\mathbf{x}^{(i)})^T & \sum_{i=1}^m \mathbf{x}^{(i)} \\ \sum_{i=1}^m \mathbf{x}^{(i)} & m \end{array} \right]^{-1} \left[ \begin{array}{c} \sum_{i=1}^m \mathbf{x}^{(i)} y^{(i)} \\ \sum_{i=1}^m y^{(i)} \end{array} \right].
 \end{aligned}$$

- General OLS Solution

- Assume that  $X \in \mathbb{R}^{m \times (n+1)}$  is a data matrix with one data case  $x^{(i)} \in \mathbb{R}^{n+1}$  per row, where  $x_0^{(i)} = 1$  and  $y \in \mathbb{R}^m$  is a column vector containing the corresponding outputs.  $w \in \mathbb{R}^{n+1}$  where  $w_0 = b$
- Use derivative:

for General OLS :  
 $b$  is included in the vector  $\vec{w}$   $\rightarrow$   $\vec{w}$  mapping all samples to a value.

$$L(\vec{w}) = \frac{1}{m} (\vec{y} - X\vec{w})^T (\vec{y} - X\vec{w})$$

$$\begin{aligned}\nabla_{\vec{w}} L(\vec{w}) &= \frac{1}{m} \nabla_{\vec{w}} (\vec{y} - X\vec{w})^T (\vec{y} - X\vec{w}) \\ &= \frac{1}{m} \nabla_{\vec{w}} \left( y^T y - \underbrace{y^T X w}_{1 \times m \times 1} - \underbrace{(X w)^T y}_{1 \times m \times 1} + (X w)^T X w \right) \\ &= \frac{1}{m} \nabla_{\vec{w}} (y^T y - 2 y^T X w + w^T X^T X w) \\ &= \frac{1}{m} \nabla_{\vec{w}} (w^T X^T X w - 2 y^T X w) \\ &= \frac{1}{m} (2 X^T X w - 2 X^T y) = \frac{2}{m} (X^T X w - X^T y) = 0 \\ w &= (X^T X)^{-1} X^T y.\end{aligned}$$

Def: ① Since the matrix multiplication is actually linear transformation.

$$\frac{\partial A\vec{x}}{\partial \vec{x}} = \frac{\partial A\vec{x}}{\partial x}. \quad (\text{the change rates are the same for each vectors})$$

If  $\vec{y}$  is  $m$  element vector,  $\vec{x}$  is  $n$  element vector.

$$\frac{\partial \vec{y}}{\partial \vec{x}} = \left[ \begin{array}{cccc} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \frac{\partial y_m}{\partial x_2} & \cdots & \frac{\partial y_m}{\partial x_n} \end{array} \right]$$

Shape:  $n \times m$

$$\text{Therefore: } \frac{\partial A\vec{x}}{\partial \vec{x}} = A^T$$

$$②. \quad \alpha = \vec{x}^T A \vec{x} \quad x(n \times 1) \quad A(n \times n)$$

$$\frac{\partial \alpha}{\partial \vec{x}} = \vec{x}^T (A + A^T)$$

$$\text{proof: } \alpha = \sum_{j=1}^n \left( \sum_{i=1}^n a_{ij} x_i \right) x_j \quad \frac{\partial \alpha}{\partial x_k} = \sum_{j=1}^n a_{kj} x_j + \sum_{i=1}^n a_{ik} x_i$$

$$\frac{\partial \alpha}{\partial \vec{x}} = \vec{x}^T A^T + \vec{x}^T A = \vec{x}^T (A^T + A) \rightarrow$$

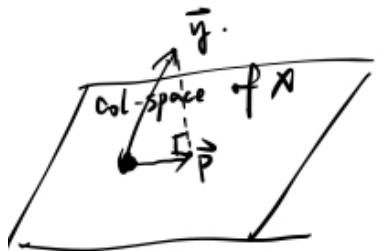
when  $A$  is symmetric.

$$\frac{\partial \alpha}{\partial \vec{x}} = 2 \vec{x}^T A$$

$$\left[ \begin{array}{c} \frac{\partial \alpha}{\partial x_1} \\ \frac{\partial \alpha}{\partial x_2} \\ \vdots \\ \frac{\partial \alpha}{\partial x_m} \end{array} \right]$$

- Use projection:

for  $\mathbf{X}\mathbf{w} = \mathbf{y}$ , there is no  $\mathbf{w}$  that fit  $\mathbf{X}$  into  $\mathbf{y}$  if  
the points in  $\mathbf{X}$  is discrete and non-linear.  
but projection  $\vec{P}$  is the best answer.  
(closest Euclidean Distance),  $\hat{\mathbf{y}} = \vec{P}$ .



$$\therefore \mathbf{X}\mathbf{w}^* = \hat{\mathbf{y}}$$

$$\hat{\mathbf{y}} - \vec{P} = (\mathbf{y} - \hat{\mathbf{y}}) \perp \mathbf{X}$$

$$\mathbf{X}^T(\mathbf{y} - \mathbf{X}\mathbf{w}^*) = 0.$$

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

- Connection to Probabilistic Models

- This same solution can be derived as the maximum conditional likelihood estimate for the parameters of a conditional Normal model
- Assumption: The output (or the observation)  $\mathbf{y}$  is from a deterministic function with additive Gaussian noise:

$$\mathbf{y} = \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon}, \text{ where } p(\boldsymbol{\epsilon}; \sigma^2) = \mathcal{N}(0, \sigma^2 \mathbf{I})$$

- This implies that

$$p(\mathbf{y} | \mathbf{X}; \mathbf{w}, \sigma^2) = \frac{1}{\sqrt{(2\pi\sigma^2)^m}} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\mathbf{w})^T(\mathbf{y} - \mathbf{X}\mathbf{w})\right)$$

- Maximizing the log conditional likelihood, we have

$$\begin{aligned}
& \max_{\mathbf{w}, \sigma^2} \log p(\mathbf{y} | \mathbf{X}; \mathbf{w}, \sigma^2) \\
&= \max_{\mathbf{w}, \sigma^2} -\frac{m}{2} \log 2\pi - \frac{m}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) \\
&= \min_{\mathbf{w}, \sigma^2} \frac{m}{2} \log \sigma^2 + \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})
\end{aligned}$$

- We note that solving for  $\mathbf{w}$  is independent of  $\sigma^2$ :

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \frac{1}{2} (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})$$

- Proof:

### 3.5 A Probabilistic View of Linear Regression

The output  $y$  is from a deterministic function with additive Gaussian noise:

$$y = Xw + \epsilon, \text{ where } p(\epsilon; \sigma^2) = \mathcal{N}(\epsilon; 0, \sigma^2 I). \quad (49)$$

Equivalently,

$$p(y|X; w, \sigma^2) = p(y - Xw; \sigma^2) = \mathcal{N}(y - Xw; 0, \sigma^2 I) = \mathcal{N}(y; Xw, \sigma^2 I). \quad (50)$$

Recall the multivariate Gaussian distribution

$$\mathcal{N}(z; \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^m |\Sigma|}} \exp\left(-\frac{1}{2}(z - \mu)^T \Sigma^{-1} (z - \mu)\right), \quad (51)$$

and we have

$$\begin{aligned}
p(y|X; w, \sigma^2) &= \frac{1}{\sqrt{(2\pi)^m |\sigma^2 I|}} \exp\left(-\frac{1}{2}(y - Xw)^T (\sigma^2 I)^{-1} (y - Xw)\right) \\
&= \frac{1}{\sqrt{(2\pi\sigma^2)^m}} \exp\left(-\frac{1}{2\sigma^2} (y - Xw)^T (y - Xw)\right).
\end{aligned} \quad (52)$$

- Normal Equations

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

- Pseudo-inverse

$$\mathbf{X}^\dagger = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

- For invertible  $X$  pseudo-inverse is just inverse
  - When  $X^T X$  is non-invertable:
    - $m < n + 1 \Rightarrow$  increase sample set and regularization
    - $m > n + 1 \Rightarrow$  singular (remove redundant features)
      - Reason:  $\text{rank}(X^T X) = \text{rank}(X)$
  - $n$  to large  $\Rightarrow$  GD
  - $m$  to large  $\Rightarrow$  SGD / mini-batch GD
  - Strengths and Limitations of OLS
    - Need at least  $m \geq n + 1$  data cases to learn a model with an  $n$  dimensional feature vector. Otherwise **inverse of  $X^T X$**  is not defined
    - Very sensitive to noise and outliers due to MSE objective function/Normally distributed residuals assumption (prefer many medium error than one large error)
    - Sensitive to co-linear features ( $x_i \approx ax_j + b$ ). Otherwise inverse of  $X^T X$  is not numerically stable
    - Computation is cubic in data dimension  $n$
- 

## Lecture06: Regularized Linear Regression: Ridge, and Lasso & Overfitting, Regularization, and Cross Validation

### Regularized Linear Regression

- Reason of regularization
  - Control the parameter space to avoid *overfitting*
  - Obtain numerically more stable solutions
  - Enforce the desired parameter space as a prior
  - What is the price that we pay? The regularization term will bias the least squares estimate (not exactly least square)
  - Probabilistic interpretation:

$$\begin{aligned}
& P(w|y, X; \sigma^2, \lambda) \propto P(y|w, X; \sigma^2) P(w; \lambda) \\
& = \frac{P(w, y, X; \sigma^2, \lambda)}{P(y|w, X; \sigma^2, \lambda) \cdot P(w, \lambda)} \\
& = \frac{P(y, X; \sigma^2)}{P(w; \lambda)} \\
& \log P(w|y, X; \sigma^2, \lambda) \propto \log P(y|w, X; \sigma^2) + \log P(w; \lambda)
\end{aligned}$$

- General framework:

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \frac{\lambda}{2} \sum_{j=1}^n |w_j|^q = \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \frac{\lambda}{2} \|\mathbf{w}\|_q^q$$

- For ridge regression,  $q = 2$ . Closed-form solution exists
- For lasso regression,  $q = 1$ . It is a quadratic programming problem with no closed-form solution in general. However, it can be solved efficiently using an algorithm called *least angle regression*. The advantage is that lasso simultaneously performs regularization and **encourages sparsity** (as a way of feature selection)
- When  $q < 1$ , the problem is non-convex, but **encourages more sparsity**
- The relationship between the **weight matrix** of regularized version and linear regression version
  - assume that:  $X^T X = I$  (each feature vector are unit orthogonal basis)
  - Linear regression:  $\min_w \frac{1}{2} \|Xw - y\|_2^2$

$$w_{LS} = (X^T X)^{-1} X^T y = X^T y$$

$$\circ \text{ Ridge regression: } \min_w \frac{1}{2} \|Xw - y\|_2^2 + \frac{\lambda}{2} \|w\|_2^2$$

$$w_{RR} = (X^T X + \lambda I)^{-1} X^T y = X^T y / (1 + \lambda) = W_{LS} / (1 + \lambda)$$

$$\circ \text{ Lasso regression: } \min_w \frac{1}{2} \|Xw - y\|_2^2 + \lambda \|w\|_1$$

$$w_{RR} = (X^T X + \lambda I)^{-1} X^T y = X^T y / (1 + \lambda) = W_{LS} / (1 + \lambda)$$

- Proof:

Ridge

$$\begin{aligned}
 & \min \left( \frac{1}{2} \|Xw - y\|_2^2 + \frac{\lambda}{2} \|w\|_2^2 \right) \\
 &= \min \left( \frac{1}{2} (Xw - y)^T (Xw - y) + \frac{\lambda}{2} w^T w \right) \quad \text{b.c. } (X^T X = X^T X^{-1} = I) \\
 &= \min \left( \frac{1}{2} (X^T (Xw - y))^T (X^T (Xw - y)) + \frac{\lambda}{2} w^T w \right) \\
 &= \min \left[ \frac{1}{2} (X^T Xw - X^T y)^T (X^T Xw - X^T y) + \frac{\lambda}{2} w^T w \right] \\
 &= \min \left[ \frac{1}{2} (w - X^T y)^T (w - X^T y) + \frac{\lambda}{2} w^T w \right] \\
 &= \min \left( \|w - X^T y\|_2^2 + \lambda \|w\|_2^2 \right) \quad \text{b.c. } X^T y = w_{LR} \\
 &= \min \left( \|w - w_{LR}\|_2^2 + \lambda \|w\|_2^2 \right) \\
 &= \min \left( \|w\|_2^2 (\lambda + 1) - 2 w^T w_{LR} + \|w_{LR}\|_2^2 \right) \\
 &= w = \left( \frac{1}{1 + \lambda} \right) w_{LR}.
 \end{aligned}$$

LASSO

$$\begin{aligned}
 & \min \left( \frac{1}{2} \|Xw - y\|_2^2 + \lambda \|w\|_1 \right) \\
 &= \min \left( \frac{1}{2} \|w\|_2^2 - w^T W_{LR} + \frac{1}{2} \|w_{LS}\|_2^2 + \lambda \|w\|_1 \right) \\
 &= \min \left( \frac{1}{2} \sum_{j=1}^n w_j^2 - \sum_{j=1}^n w_j (w_{LS})_j + \lambda \sum_{j=1}^n |w_j| + \frac{1}{2} \|w_{LS}\|_2^2 \right) \\
 &\quad \text{for each } j: f(w_j) = \frac{1}{2} w_j^2 - w_j (w_{LS})_j + \lambda |w_j| + \frac{1}{2} (w_{LS})_j^2 \\
 &\quad \text{when } w_j (w_{LS})_j < 0 \\
 &\quad -w_j (w_{LS})_j > 0 \Rightarrow -(-w_j) (w_{LS})_j \\
 &\quad \text{Therefore, } -w_j \text{ is a better solution than } w_j. \\
 &\quad \text{We only need to find the optimal under the constraint } w_j (w_{LS})_j \geq 0
 \end{aligned}$$

$\textcircled{1}$  when  $(w_{LS})_j > 0 \quad w_j \geq 0 \quad \text{if } (w_{LS})_j - \lambda < 0$

$$\begin{aligned}
 f(w_j) &= \frac{1}{2} w_j^2 + (\lambda - (w_{LS})_j) w_j \\
 w_j^* &= -\frac{\lambda - (w_{LS})_j}{2 \times \frac{1}{2}} = -(\lambda - (w_{LS})_j) \\
 &= (w_{LS})_j - \lambda \quad (\text{if } (w_{LS})_j - \lambda \leq 0)
 \end{aligned}$$

$\textcircled{2}$  when  $(w_{LS})_j < 0 \quad w_j \leq 0$

$$\begin{aligned}
 f(w_j) &= \frac{1}{2} w_j^2 - (\lambda + (w_{LS})_j) w_j \\
 w_j^* &= \frac{\lambda + (w_{LS})_j}{2 \times \frac{1}{2}} = \lambda + (w_{LS})_j \quad (\text{when } \lambda + (w_{LS})_j \leq 0) \\
 w_j^* &= \min (\lambda + (w_{LS})_j, 0)
 \end{aligned}$$

From  $\textcircled{1}$  &  $\textcircled{2}$ .

$w_j^* = \text{sign}(w_{LS}) \max(0, |w_{LS}| - \lambda)$

$w_j^* = 0 \quad \text{if } (w_{LS})_j = 0 \quad \text{is the best}$

$\text{if } \lambda + (w_{LS})_j > 0.$

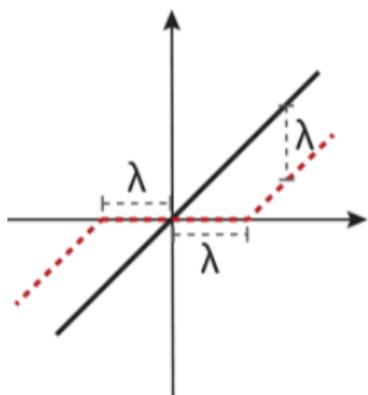
$\text{if } w_j^* < 0 \Rightarrow -(\lambda + (w_{LS})_j) w_j^* > 0.$

$f(w_j) > 0$

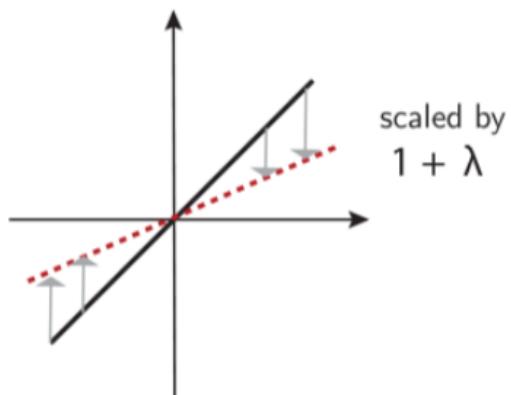
$w_j^* = 0 \Rightarrow f(w_j) = 0$

- Geometrical meaning (with x-axis  $W_{LR}$ )

LASSO



Ridge



- Lasso vs Ridge (2D case)

- Property of the minimum point

$$l_r = \sum_{i=1}^n (y - x_i^{(i)} w_1 - x_i^{(i)} w_2 - b) + \lambda (|w_1|^r + |w_2|^r) = C$$

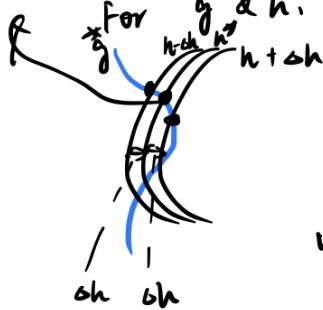
if the optimal of  $l_r$  is found to be  $C$   
and at that moment

$$\begin{cases} \sum_{i=1}^n (y - x_i^{(i)} w_1^* - x_i^{(i)} w_2^* - b) = C_1, \\ \lambda (|w_1^*|^r + |w_2^*|^r) = C_2 \end{cases}$$

on  $w_1-w_2$  plane, they will be two curves,  
which meet at the point  $(w_1^*, w_2^*)$ .

$(w_1^*, w_2^*)$  is the only common point

e.g.  $f(w_1, w_2) = g(w_1, w_2) + h(w_1, w_2)$ .  
 $(w_1^*, w_2^*)$  for  $g$  &  $h$ , if they are not intersect at only 1 point



for minimum value.  
when  $ah \rightarrow 0$ ,  $h+ah$  and  $h-ah$  will have similar shape with  $h^*$ .

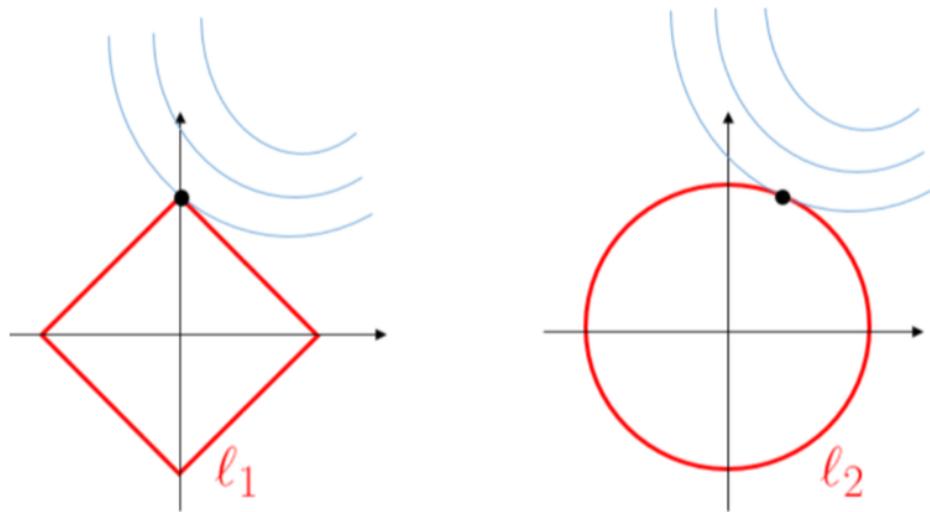
when  $ah > 0$ ,  $\frac{ah}{h+ah} < h^*$        $\frac{ah}{h-ah} < h^*$

$$f = \tilde{g} + h^* - ah < f^* \quad f = \tilde{g}^* + h^* + ah < f^*.$$

Therefore, there is only one intersection.

- Graph of Lasso regression and Ridge regression:

## Why $\ell_1$ leads to sparsity: 2D case



- it's easier to have the intersection on the corner for  $\ell_1$  case
- one of  $w_j$  will be zero, which means the feature is dropped => sparsity
- Strengths and Limitations of Ridge and Lasso
  - Need at least  $m \geq n + 1$  data cases to learn a model with an  $n + 1$  dimensional feature vector (similar as linear regression)
  - Work when  $X^T X$  is close to singular => select feature
  - MSE objective function is still sensitive to noise points
  - Do not solve **bias problem**
  - Computation for ridge is still cubic in data dimension  $n$  with additional cost to determine regularization parameters. Computation for lasso is similar

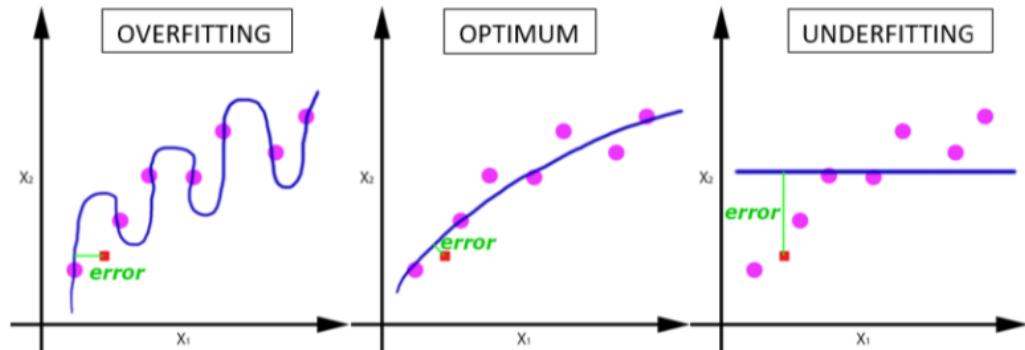
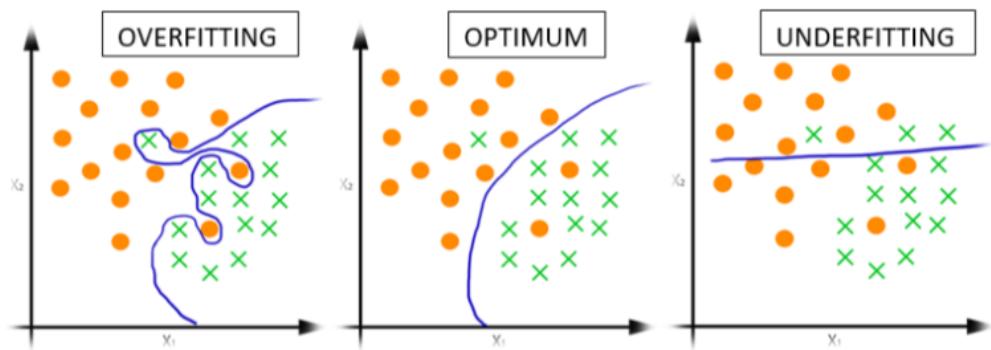
## Review

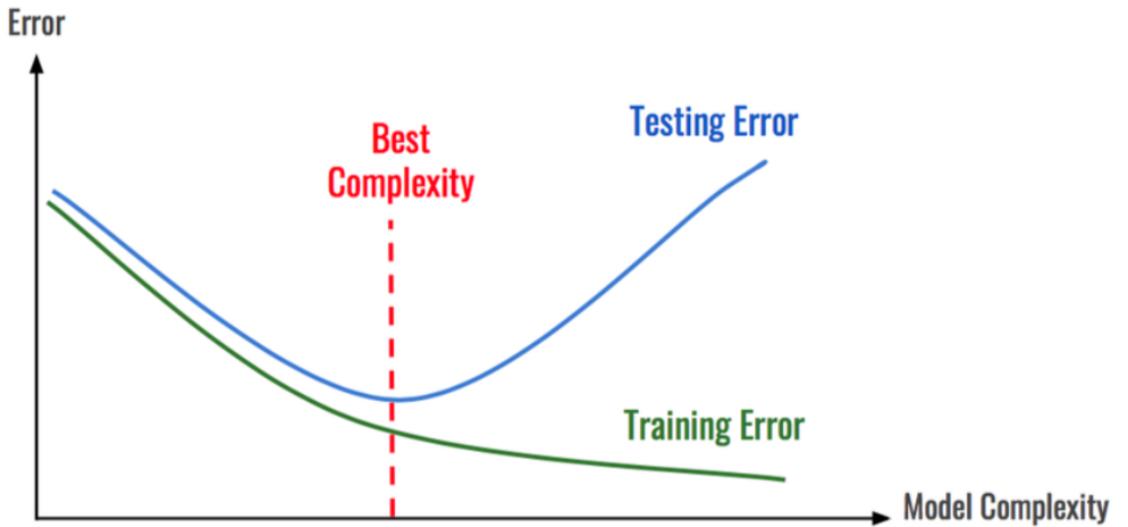
- We have learnt 4 classifiers in depth: KNN, Naive Bayes, LDA and Logistic Regression
  - KNN: Results sensitive to  $k$
  - Logistic Regression:  $p(y|x; w)$  sensitive to the magnitude of  $w$

## Capacity and Generalization

- Model capacity:
  - **Deterministic classifiers** that can represent **more complex decision boundaries** (in log form) are said to have **higher capacity** than classifiers that can only represent simpler boundaries
  - **Probabilistic classifiers** that can represent **more complex sets of conditionals**  $p(y|x)$  are said to have higher *capacity* than probabilistic classifiers that can only represent simpler sets of conditionals
  - Rank of the learned classifiers: KNN > Naive Bayes, Logistic Regression  $\approx$  LDA
- Reason of not always using the classifier with the highest possible capacity for every problem
  - This would minimize the error on the **training data**. However, what we really care about for prediction problems is *generalization*

- **Generalization:** The ability of a trained classifier to achieve an error rate on *future, unseen examples* (the generalization error rate) that is comparable to the training error rate
- **Capacity Control:** To achieve optimal generalization performance for a given training set, we often need to control model capacity carefully
- **Overfitting vs underfitting**
  - **Overfitting:** The generalization error for a classifier is much worse than the training error. This usually results from choosing a classifier with too much capacity so that it models the noise in the training data
  - **Underfitting:** Occurs when the capacity of the classifier is too low to capture the actual structure in the training data, leading to both high training error and high generalization error
  - Examples(error on the down-left corner graph):





- Bias-Variance Trade-Off (both of them are negative things)
  - Bias: A classifier is said to have low *bias* if the true decision boundary or conditionals  $p(y|x)$  can be **approximated closely by the model**
  - Variance: A classifier is said to have low *variance* if the decision boundary or conditionals  $p(y|x)$  it constructs are stable with respect to small changes in the training data
  - Bias-Variance Dilemma: To achieve **low generalization error**, we need classifiers that are low-bias and low-variance, but this isn't always possible
  - Bias-Variance and Capacity: On **complex data**, models with **low capacity have low variance**, but **high bias**; while models with **high capacity** have **low bias**, but **high variance**

## Hyperparameters

- Why need to control hyperparameters
  - In order to control the capacity of a classifier, it needs to have capacity control parameters
  - Because capacity control parameters can not be chosen based on training error, they are often called hyperparameters
  - For KNN:  $k$ , for logistic regression:  $\alpha$  and  $\lambda$
- Capacity, Smoothness and Regularization

for LR , if  $w$  is large (magnitude)

$$P(y|x;w) = \frac{1}{1+e^{-w^T x}} \geq 0.5 \Leftrightarrow y=1$$

$$w^T x > 0 \Rightarrow y=1$$

$$w^T x < 0 \Rightarrow y=0$$

$$(x_i - \bar{x})w_i = x_i w_i - \bar{x} w_i$$

large  $\Rightarrow$  change  $x$ , validate  $\Rightarrow$  different result.

- **Capacity and Smoothness:** In the case of probabilistic classifiers like naive Bayes, LDA, and logistic regression, we can think of capacity in terms of the smoothness of  $p(y|x)$
- **Regularization:** We can control the smoothness of  $p(y|x)$  using a technique called *regularization* that penalizes parameters that result in overly complex  $p(y|x)$  during learning
- **Laplace Smoothing:** In the case of Naive Bayes, we introduced Laplace smoothing for a categorical distribution

$$\theta_{x,j|y} = \frac{\sum_{i=1}^m \mathbb{I}[y^{(i)} = y \cap x_j^{(i)} = x] + \alpha}{\sum_{i=1}^m \mathbb{I}[y^{(i)} = y] + \alpha |\mathcal{X}|}$$

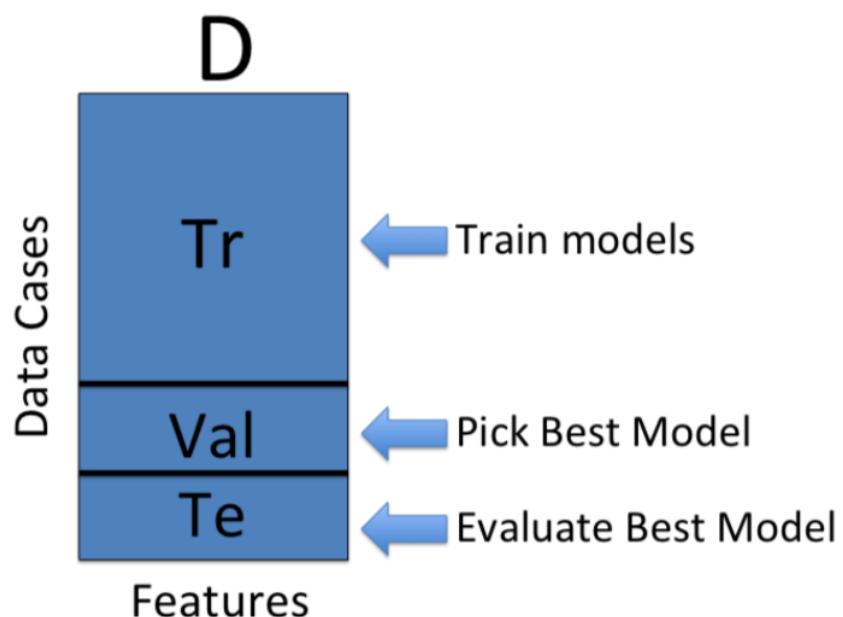
- As the regularization hyperparameter  $\alpha$  increases, our estimate of the distribution smooths out toward being uniform
- In the case of logistic regression, the smoothness of  $p(y|x; \theta)$  is determined by the magnitude of  $\theta$ 
  - Similar in linear regression, we can control the magnitude of  $\theta$  by introducing a regularization term into the conditional log likelihood learning criteria that penalizes the  $l_p$ -norm of  $\theta$

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} - \sum_{i=1}^m \log p(y^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta}) + \lambda \|\boldsymbol{\theta}\|_p^p$$

- As the regularization hyperparameter  $\lambda$  increases, the learned  $p(y|x; \theta)$  will smooth out

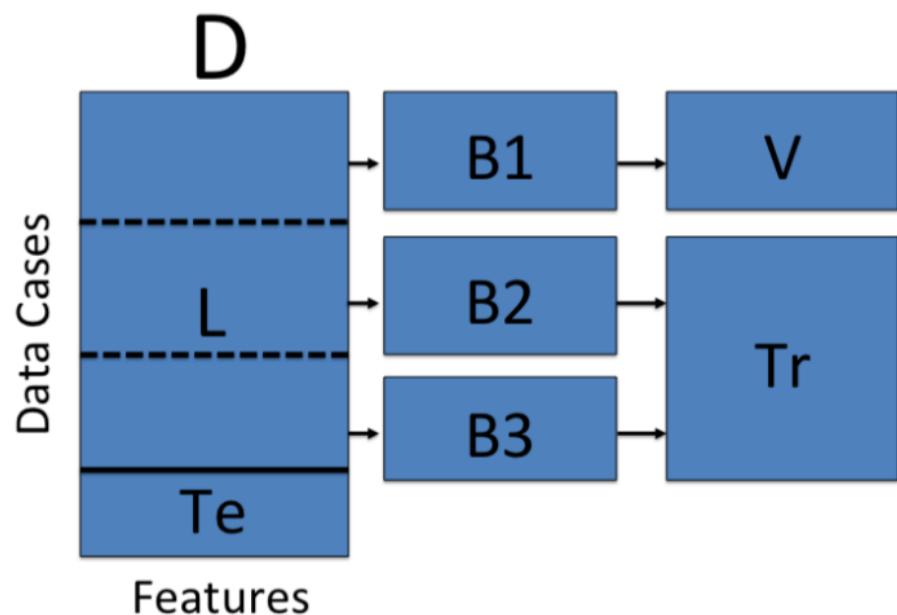
# Model Selection and Evaluation

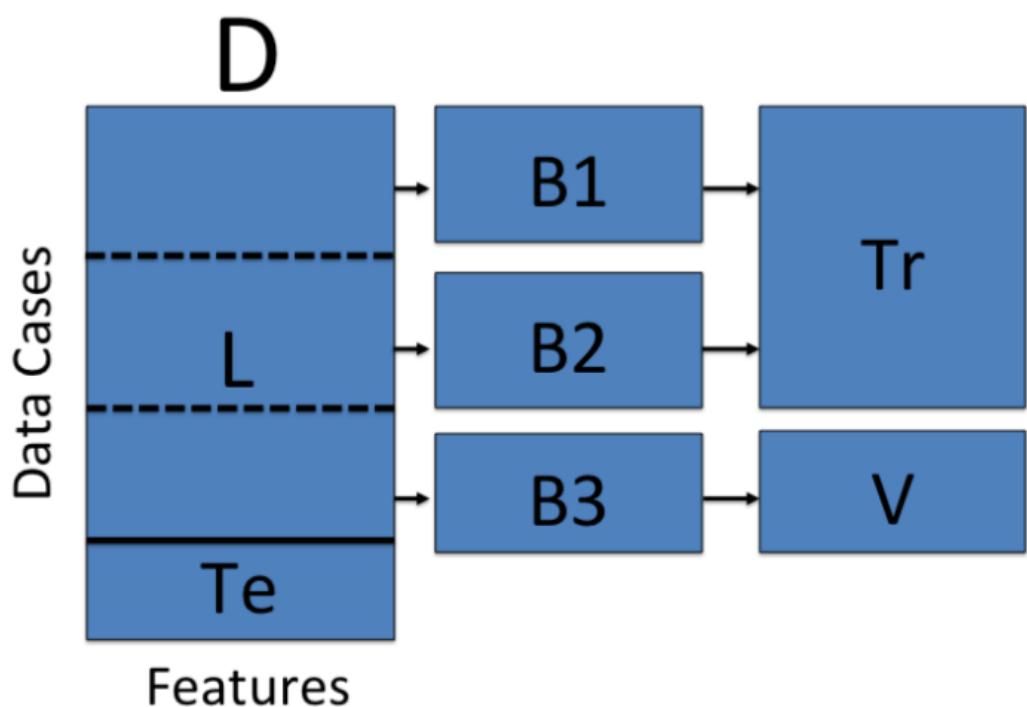
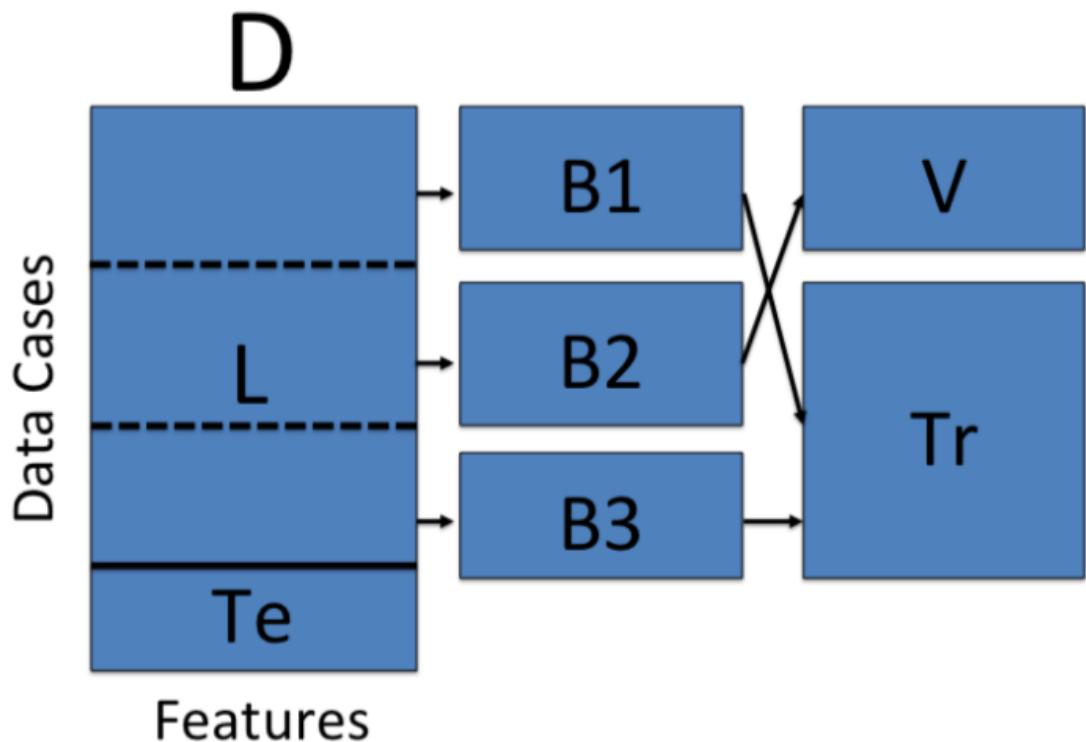
- Principles:
  - We've identified the parameters that **control the capacity** of our models, we need a way to choose optimal values for these parameters
  - In addition, we will want an estimate of the **generalization error** that the selected parameters achieve
  - To obtain valid results, we need to use appropriate methodology and construct learning experiments carefully
  - **Guiding Principle:** Data used to estimate generalization error(**testing set**) can not be used for any other purpose, *i. e.*, model training, hyperparameter selection, feature selection, etc. Otherwise, the results of the evaluation will be *biased*
- Methods:
  - Recipe 1: Train, Validation, and Test
    - Given a data set  $D$ , we randomly partition the data cases into a **training set (Tr)**, a **validation set (Val)**, and a **test set (Te)**. Typical splits are 60/20/20, 90/5/5, etc.
    - Training
      - Models  $\{f(i)\}$  are learned on Tr for each **one** choice of hyperparameters  $\theta(i)$  (Choose weight matrix in training set)
    - Validation
      - The validation error  $Err_v^{(i)}$  of each model  $f^{(i)}$  is evaluated on Val  $v$
      - The hyperparameters  $\theta^*$  with the lowest validation error are selected and the classifier is re-trained using these hyperparameters on **Tr + Val**, yielding a final model  $f^*$  (Choose hyperparameters in training & validation set)
    - Testing
      - Generalization performance is estimated by evaluating error/accuracy of  $f^*$  on the test data Te (Evaluate the final model)
  - Example: (Notice that the data cases needs to be randomly shuffled before partitioning  $D$ )



- Recipe 2: Cross-Validation and Test

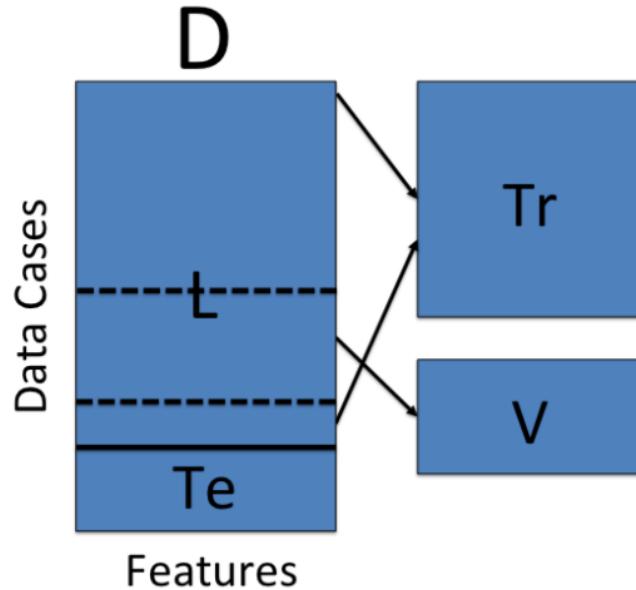
- Randomly partition  $D$  into a learning set  $L$  and a test set  $Te$ , typically 50/50, 80/20, etc.
- We next **randomly partition**  $L$  into a set of  $K$  blocks  $B_1, \dots, B_K$
- For each cross-validation fold  $k = 1, \dots, K$ :
  - Let  $Val = B_k$  and  $Tr = L/B_k$  (the remaining  $K - 1$  blocks)
  - Learn  $f^{(i,k)}$  on  $Tr$  for each choice of hyperparameters  $\theta^{(i,k)}$
  - Compute  $Err_v^{(i,k)}$  of  $f^{(i,k)}$  on  $Val$
- Select hyperparameters  $\theta^*$  minimizing  $\frac{1}{K} \sum_{k=1}^K Err_v^{(i,k)}$  and re-train model on  $L$  using these hyperparameters, yielding the final model  $f^*$
- Estimate generalization performance by evaluating error/accuracy of  $f^*$  on  $Te$
- Example: Three Fold Cross-Validation and Test (Notice that the data cases needs to be randomly shuffled before partitioning  $D$  into  $L$  and  $Te$ )





- Recipe 3: Random Resampling Validation and Test (Notice that the data cases needs to be randomly shuffled before partitioning  $D$  into  $L$  and  $Te$ )
  - Randomly partition the data cases into a learning set  $L$  and a test set  $Te$ , typically 50/50, 80/20, etc.
  - For sample  $k = 1, \dots, K$ :
    - Randomly partition  $L$  into  $Tr$  and  $Val$ , again 50/50, 80/20, etc.

- Learn  $f^{(i,k)}$  on Tr for each choice of hyperparameters  $\theta^{(i,k)}$
- Compute  $Err^{(i,k)}$  of  $f^{(i,k)}$  on Val  $v$
- Select hyperparameters  $\theta^*$  minimizing  $\frac{1}{K} \sum_{k=1}^K Err_v^{(i,k)}$  and re-train model on  $L$  using these hyperparameters, yielding the final model  $f^*$
- Estimate generalization performance by evaluating error/accuracy of  $f^*$  on Te
- Example:



- Trade-Offs:
  - In cases where the data has a benchmark split into a training set and a test set, we can use Recipes 1-3 by preserving the given test set and splitting the given training set into train and validation sets as needed
  - Choosing larger  $K$  in cross-validation will reduce bias. Choosing larger  $K$  in random re-sampling validation will reduce variance and bias. However, both increase computational costs.  $K = 3, 5, 10$  are common choices for cross-validation. ,  $K = m$ , also known as **leave-one-out cross validation** is also popular when feasible

## Feature Selection

- Feature Selection: As a special and important case of model selection, feature selection selects a useful subset from all the features
- Reason of feature selection:
  - Some algorithms (e.g. normal equation) scale (computationally) poorly with increased feature dimension
  - Irrelevant features can confuse some algorithms
  - Redundant features may adversely affect regularization
  - Reduces data set and resulting model size
- Methods
  - Wrapper methods (keep the learning algorithm in the loop)
    - Requires repeated runs of the learning algorithm with different sets of features

- Can be computationally expensive
- Filter feature selection methods
  - Use some ranking criteria to rank features
  - Select the top ranking features
- Wrapper Methods
  - Forward search
    - Start with no features
    - Greedily include the most relevant feature
    - Stop when selected the desired number of features
    - Steps:
      - Let  $F = \{\}$
      - While not selected desired number of features
      - For each unused feature  $f$ 
        - Estimate model's error on feature set  $F \cup f$
        - Add  $f$  with lowest generalization error to  $F$
  - Backward search (Maybe computational expensive)
    - Start with all the features
    - Greedily remove the least relevant feature
    - Stop when selected the desired number of features
    - Steps:
      - Let  $F = \{all\ features\}$
      - While not reduced to desired number of features
      - For each unused feature  $f \in F$ 
        - Estimate model's error on feature set  $F \setminus f$  (using cross-validation)
        - Remove  $f$  with lowest generalization error from  $F$
  - Inclusion/Removal criteria uses cross-validation (*i.e.*, train model using the current subset of features and estimate the generalization error)
- Filter Feature Selection Methods (outdated => deep learning nowadays)
  - Uses heuristics but are much faster than wrapper methods
  - Correlation criterion: Rank features according to their correlation with the labels:

$$\text{Corr}(\mathbf{x}_j, \mathbf{y}) = \frac{\sum_{i=1}^m (x_j^{(i)} - \mu_{\mathbf{x}_j})(y^{(i)} - \mu_{\mathbf{y}})}{\sqrt{\sum_{i=1}^m (x_j^{(i)} - \mu_{\mathbf{x}_j})^2} \sqrt{\sum_{i=1}^m (y^{(i)} - \mu_{\mathbf{y}})^2}}$$

- Mutual information criterion:

$$\text{MI}(x_j; y) = \sum_{x_j \in \mathcal{X}_j} \sum_{y \in \mathcal{Y}} p(x_j, y) \log \frac{p(x_j, y)}{p(x_j)p(y)},$$

- where the probabilities  $p(x_j, y)$ ,  $p(x_j)$  and  $p(y)$  can be estimated according to their empirical distributions on the training set
- More on Mutual Information
  - Mutual information is a fundamental concept in *information theory*, measuring the mutual dependence between two variables
  - It quantifies the “amount of information” obtained about one random variable through observing the other random variable
  - As a function of  $p(x, y) = p(x)p(y|x)$ ,  $I(x; y)$  is concave in  $p(x)$  for fixed  $p(y|x)$ , and is convex in  $p(y|x)$  for fixed  $p(x)$
  - Maximizing  $I(x; y)$  w.r.t.  $p(x)$  leads to *noisy-channel coding theorem*, establishing the maximum communication rate through a noisy channel
  - Minimizing  $I(x; y)$  w.r.t.  $p(y|x)$  leads to *rate-distortion theorem*, laying the theoretical foundations for lossy data compression

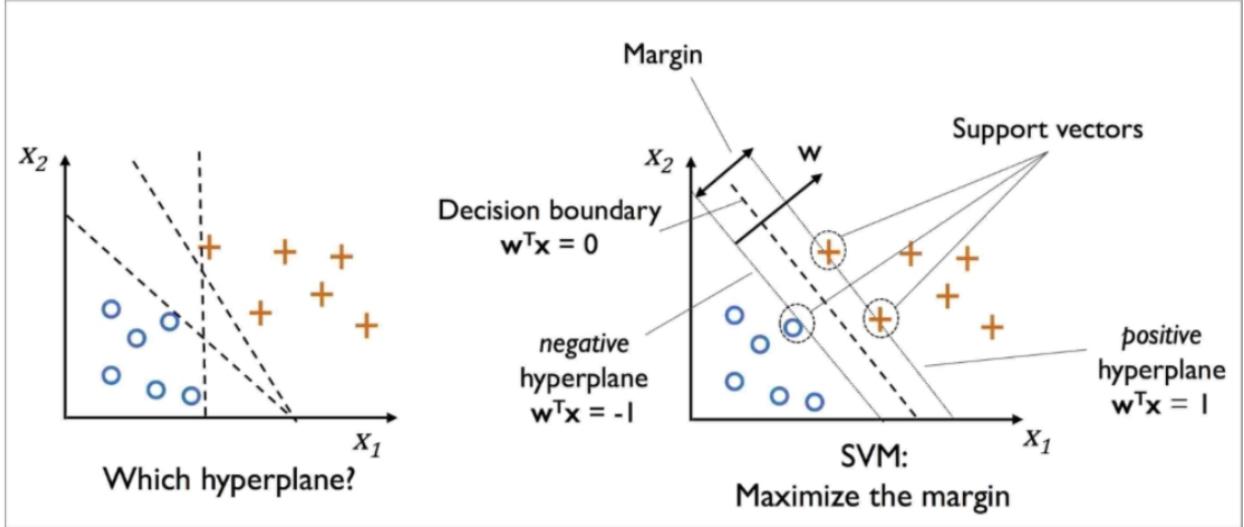
## Lecture07 Support Vector Machine

Target: Increase the capacity of linear classifiers so they can produce non-linear classification boundaries

### Support Vector Machine (SVM)

- Decision Boundary
  - A binary SVM is a discriminative classifier that takes labels in the set  $\{-1, 1\}$
  - The decision function has the form:
 
$$f_{SVM}(x) = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$$

$$\mathbf{w}^T \mathbf{x} + b = 0 \text{ (on the decision boundary)}$$
  - It’s easy to show that the decision boundary for logistic regression can be written in exactly the same way
  - Question: If logistic regression and SVMs have the same form for their decision boundaries, how do they differ
    - Logistic regression minimize the log probability of miss-classification
    - Logistic regression maximize the probability on the correct label (i.e. if the probability is more than 0.5, it can be correctly classified. But it tends to increase the probability during training)
- Maximum Margin Principle (**notice the axis**)



- Define the *margin* of a linear classifier as the width that the **boundary could be increased before hitting a data point**
- The simplest kind of SVM is a linear classifier with the maximum margin
- Reason of MAX margin:
  - If we've made a small error in the location of the decision boundary, this gives us least chance of causing a misclassification
  - There's solid theory (using VC dimension) that provides indirect support for this proposition
- Compute the Margin
  - The distance between a point to a line:
$$d^{(i)} = \frac{|\mathbf{w}^T \mathbf{x} + b|}{\|\mathbf{w}\|_2} = \frac{y^{(i)}(\mathbf{w}^T \mathbf{x} + b)}{\|\mathbf{w}\|_2}$$
  - Proof:

$$ax + by + cz + d = 0.$$

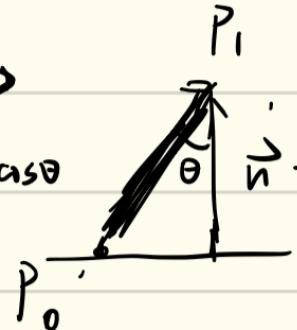
$$\text{set } \vec{b} = \langle x_1 - x_0, y_1 - y_0, z_1 - z_0 \rangle$$

$$P = |\text{comp}_{\vec{n}} \vec{b}| = \frac{|\vec{n} \cdot \vec{b}|}{|\vec{n}|} = |\vec{b}| \cos \theta$$

$$= \frac{|a(x_1 - x_0) + b(y_1 - y_0) + c(z_1 - z_0)|}{\sqrt{a^2 + b^2 + c^2}}$$

$$= \frac{|(ax_1 + by_1 + cz_1) - (ax_0 + by_0 + cz_0)|}{\sqrt{a^2 + b^2 + c^2}}$$

since  $P_0$  lies in  $ax + by + cz + d = 0$ .



$$\left| \frac{ax_1 + by_1 + cz_1 + d}{\sqrt{a^2 + b^2 + c^2}} \right|$$

- The  $y^{(i)}$  can be added because the value can be only found in the set  $\{-1, 1\}$
- The margin of  $\mathbf{w}^T \mathbf{x} + b = 0$  with respect to a training set  $D = \{(\mathbf{x}^{(i)}, y^{(i)}), i = 1, \dots, m\}$ , where  $\|\mathbf{w}\|_2 = \sqrt{\sum_{j=1}^n w_j^2}$

$$d(\text{margin}) = \min d^{(i)} = \min_{(\mathbf{x}, y) \in D} \frac{y^{(i)} (\mathbf{w}^T \mathbf{x} + b)}{\|\mathbf{w}\|_2}$$

- Maximize the Margin

- The maximum margin solution is found by solving

$$\max_{\mathbf{w}, b} \left( \frac{1}{\|\mathbf{w}\|_2} \min_{(\mathbf{x}, y) \in D} y^{(i)} (\mathbf{w}^T \mathbf{x} + b) \right)$$

- Note that if we make the rescaling  $\mathbf{w} \rightarrow \gamma \mathbf{w}$  and  $b \rightarrow \gamma b$ , then the objective is unchanged
- We can use this freedom to set:  $\min_{(\mathbf{x}, y) \in D} y (\mathbf{w}^T \mathbf{x} + b) = 1$

$$Ax + By + C = 0 \iff 2Ax + 2By + 2C = 0$$

- Or equivalently:  $y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1, \forall i = 1, \dots, m$

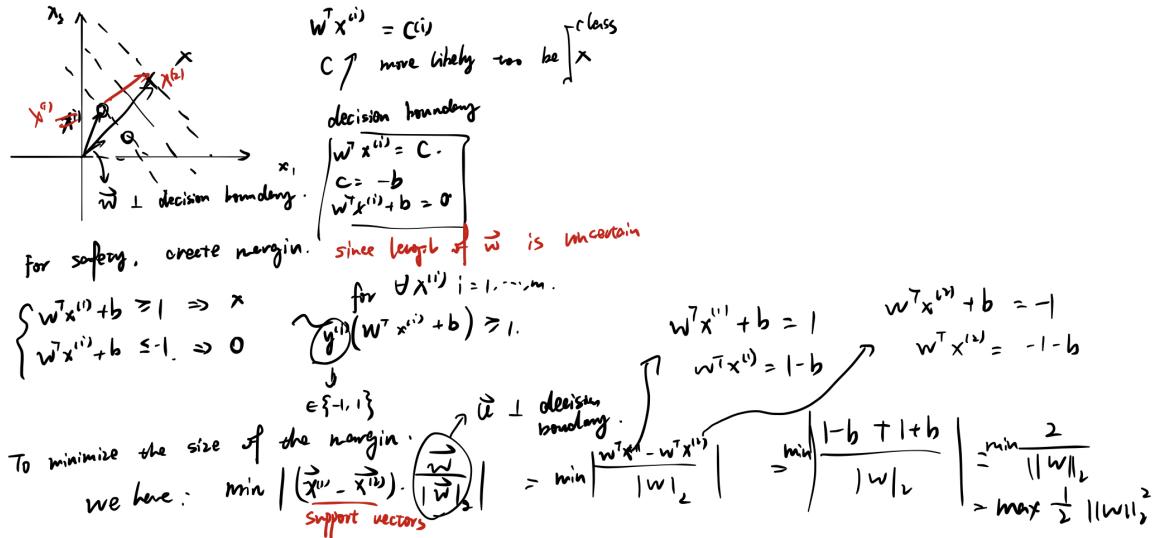
- Primal Form, Hard-Margin of SVM

$$\max_{\mathbf{w}, b} \left( \frac{1}{\|\mathbf{w}\|_2} \right) \iff \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2 \quad (\text{for convenience})$$

$$\text{subject to } y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1, \forall i = 1, \dots, m$$

- Quadratic programming problem with a quadratic objective function and linear inequality constraint
- Well studied problem in operations research

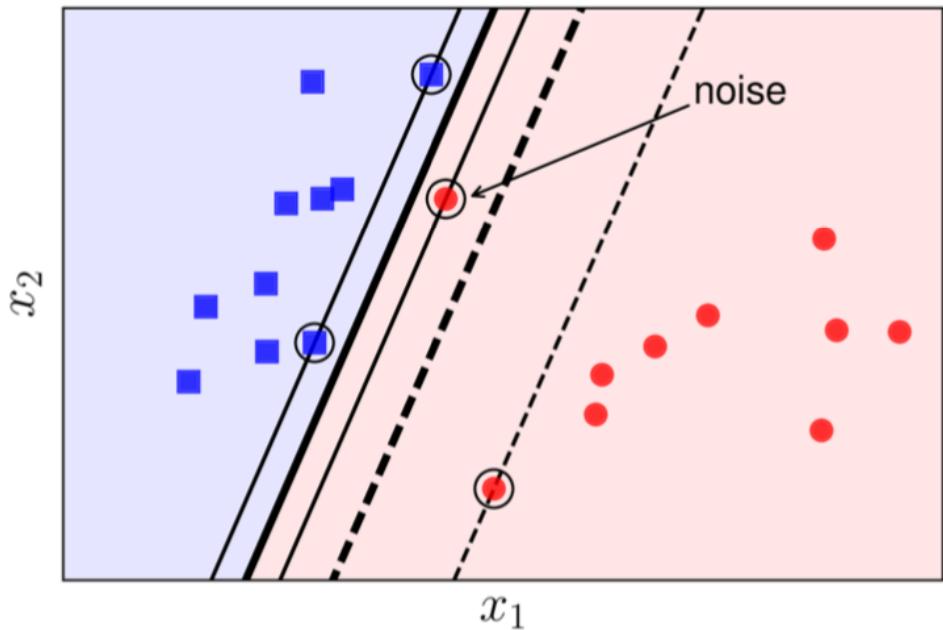
- Another point of view: (min and max in the last line should be reversed)



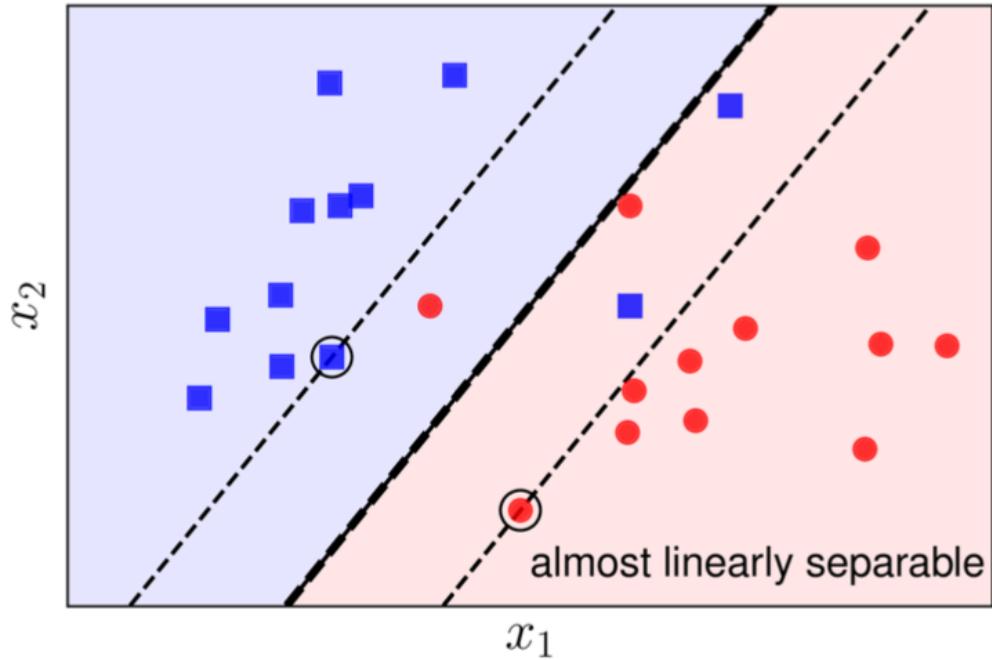
- Primal Form, Soft-Margin

- Hard margin is not enough:

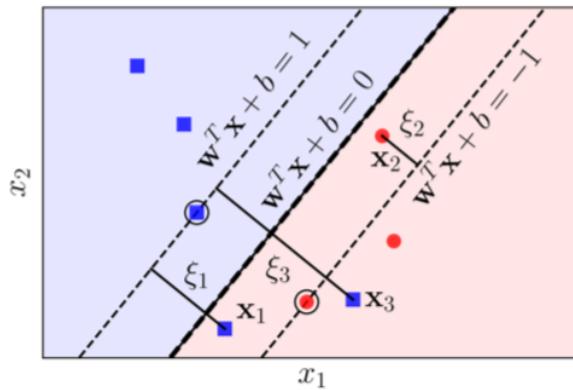
- When there is a noise point to narrow down the margin



- when the data is not linear-discriminable: notice that in this time, margin is nearly zero



- We need to add some level of tolerances to the model:



$$\begin{aligned}
 & \min_{\mathbf{w}, b} \quad \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^m \xi_i \\
 & \text{subject to} \quad y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 - \xi_i, \quad i = 1, \dots, m \\
 & \quad \xi_i \geq 0, \quad i = 1, \dots, m
 \end{aligned}$$

- In other word, for normal case,  $y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1$ , the penalty of misclassification is zero, while on the other hand, there is a penalty with the value  $1 - y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b)$
- The level of tolerance can be found on the second equation:
  - Originally,  $y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1$  means that, all pairs of the points should obey this rule, therefore, **the “distance” between two boundaries generated by the support vector should be at least 2**

- Now,  $y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 - \xi_i$ ,  $\xi_i > 0$  means that, for any one of the points, there is a **tailored**  $\xi_i$  to make sure that the point always lies in the safe constraint. Once  $y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b) < 1$ , the difference between 1 and  $y^{(i)}(\mathbf{w}^T \mathbf{x}^{(i)} + b)$  will be counted as penalty.
- This approach makes sure that all points are not forced to have a closest support vector on one or two sides of the decision boundary, but the mistake caused by it will be countered as part of the loss function

- Hinge Loss

- In the case of primal-form SVM with soft-margin, it is equivalent to minimizing the function:

$$C \sum_{i=1}^m \max(0, 1 - y^{(i)} \cdot g(\mathbf{x}^{(i)})) + \frac{1}{2} \|\mathbf{w}\|_2^2$$

- where  $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$  and  $C$  is the hyperparameter that controls the amount of regularization
- The function  $l_h(y, g(\mathbf{x})) = \max(0, 1 - y \cdot g(\mathbf{x}))$  is called the *hinge loss*

- Logistic Loss

- In the case of logistic regression with  $l_2$  regularization, we select the model parameters by minimizing the function

$$C \sum_{i=1}^m -\log p(y^{(i)} | \mathbf{x}^{(i)}) + \frac{1}{2} \|\mathbf{w}\|_2^2$$

- Under the assumption that the labels take the values  $\{-1, 1\}$ :

$$\begin{aligned} P(y^{(i)}=1 | \mathbf{x}^{(i)}) &= \frac{1}{1 + e^{-y^{(i)} g(\mathbf{x}^{(i)})}} & g(\mathbf{x}) &= \mathbf{w}^T \mathbf{x} + b \\ P(y^{(i)}=-1 | \mathbf{x}^{(i)}) &= \frac{1}{1 + e^{y^{(i)} g(\mathbf{x}^{(i)})}} & g(\mathbf{x}) &= \mathbf{w}^T \mathbf{x} + b \\ -\log p(y^{(i)} | \mathbf{x}^{(i)}) &= -\log \left( 1 + e^{-y^{(i)} g(\mathbf{x}^{(i)})} \right) \end{aligned} \Rightarrow P(y^{(i)} | \mathbf{x}^{(i)}) = \frac{1}{1 + e^{-y^{(i)} g(\mathbf{x}^{(i)})}}$$

$$C \sum_{i=1}^m \log(1 + \exp(-y^{(i)} \cdot g(\mathbf{x}^{(i)}))) + \frac{1}{2} \|\mathbf{w}\|_2^2$$

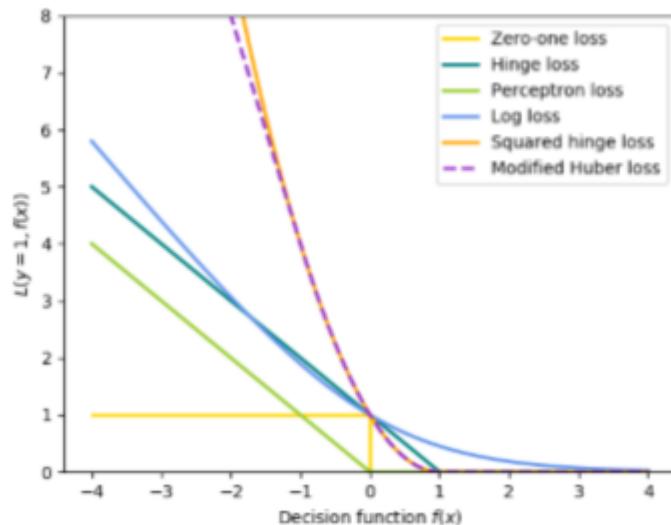
- Where the function  $l_l(y, g(\mathbf{x})) = \log(1 + \exp(-y \cdot g(\mathbf{x})))$  is called the *logistic loss function*

- Zero-One Loss

- Both the logistic loss and the hinge loss are convex upper bounds on the zero-one loss:

$$l_{01}(y, g(\mathbf{x})) = \mathbb{I}[y \neq \text{sign}(g(\mathbf{x}))]$$

- Average of 01-loss is exactly the classification error rate
- This is the loss function we'd like to minimize, but this generally isn't computationally feasible, thus the need for surrogate loss functions
- Comparison (x-axis is the decision function)



## Lagrange Duality

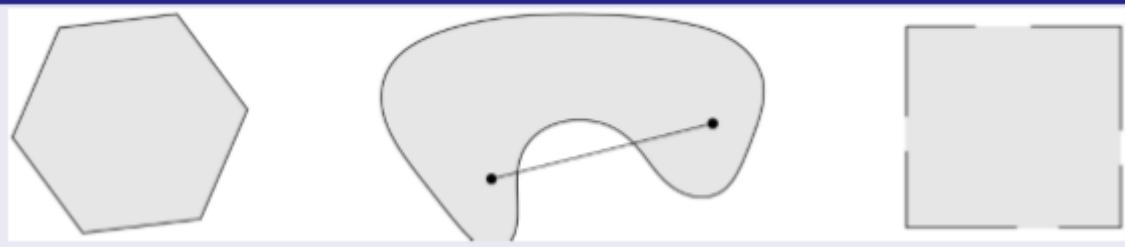
- Lagrange Duality: In [mathematical optimization](#) theory, **duality** or the **duality principle** is the principle that [optimization problems](#) may be viewed from either of two perspectives, the **primal problem** or the **dual problem**. The solution to the dual problem provides a lower bound to the solution of the primal (minimization) problem. However in general the optimal values of the primal and dual problems need not be equal. Their difference is called the [duality gap](#). For [convex optimization](#) problems, the duality gap is zero under a [constraint qualification](#) condition.
- Convex set

### Definition: Convex Set

A convex set contains line segment between any two points in the set

$$\mathbf{w}, \tilde{\mathbf{w}} \in W, 0 \leq \theta \leq 1 \implies \theta\mathbf{w} + (1 - \theta)\tilde{\mathbf{w}} \in W$$

### Example



- Convex function

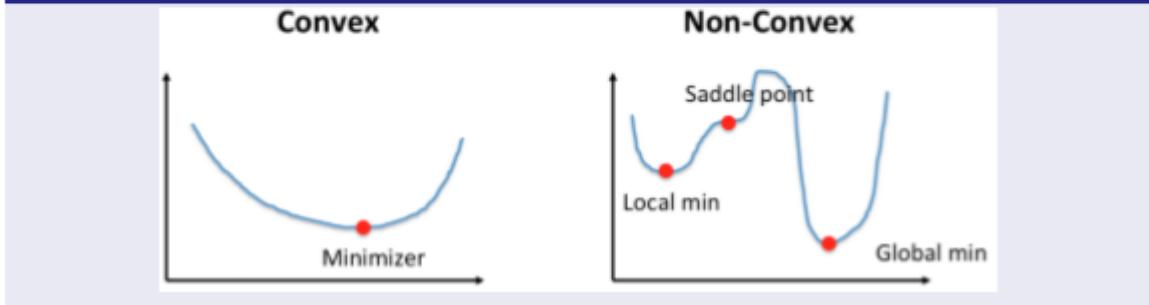
## Definition: Convex Function

$f : \mathbb{R}^n \mapsto \mathbb{R}$  is convex if the domain of  $f$ ,  $\text{dom}(f)$ , is a convex set and

$$f(\theta\mathbf{w} + (1 - \theta)\tilde{\mathbf{w}}) \leq \theta f(\mathbf{w}) + (1 - \theta)f(\tilde{\mathbf{w}})$$

for all  $\mathbf{w}, \tilde{\mathbf{w}} \in \text{dom}(f)$ ,  $0 \leq \theta \leq 1$ .

## Example



- Standard form of optimization (minimize for instance) problems

$$\begin{aligned} & \min_{\mathbf{w}} f_0(\mathbf{w}) \\ & \text{subject to } f_i(\mathbf{w}) \leq 0, i = 1, \dots, r \\ & h_i(\mathbf{w}) = 0, i = 1, \dots, s \end{aligned}$$

- $\mathbf{w} \in \mathbb{R}^n$  is the optimization variable
- $f_0 : \mathbb{R}^n \mapsto \mathbb{R}$  is the objective or cost function
- $f_i : \mathbb{R}^n \mapsto \mathbb{R}, i = 1, \dots, r$  are the inequality constraint function
- $h_i : \mathbb{R}^n \mapsto \mathbb{R}, i = 1, \dots, s$  is the equality constraint function
- If  $f_0$  and a set of  $f_i$  are convex functions, and  $h_i$ 's are affine functions, it becomes a **convex optimization problem**
  - Affine function:  $f(\mathbf{x}) = A\mathbf{x} + \mathbf{b}, A \in \mathbb{R}^{s \times n}, \mathbf{x} \in \mathbb{R}^n, \mathbf{b} \in \mathbb{R}^s$
- Lagrangian:  $L : \mathbb{R}^n \times \mathbb{R}^r \times \mathbb{R}^s \mapsto \mathbb{R}$ , with  $\text{dom}(L) : W \times \mathbb{R}^r \times \mathbb{R}^s$  (fixed weight matrix)

$$L(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\nu}) = f_0(\mathbf{w}) + \sum_{i=1}^r \lambda_i f_i(\mathbf{w}) + \sum_{i=1}^s \nu_i h_i(\mathbf{w})$$

- It's a weighted sum of objective and constraint function
- $\lambda_i$  is Lagrange multiplier associated with  $f_i(x) \leq 0$
- $\nu_i$  is Lagrange multiplier associated with  $h_i(x) = 0$
- Lagrange dual function:  $g : \mathbb{R}^r \times \mathbb{R}^s \mapsto \mathbb{R}$ , ( $\inf$  means  $\min$ )

$$g(\boldsymbol{\lambda}, \boldsymbol{\nu}) = \inf_{\mathbf{w} \in \mathbf{W}} (L(\mathbf{w}, \boldsymbol{\lambda}, \boldsymbol{\nu})) = \inf_{\mathbf{w} \in \mathbf{W}} (f_0(\mathbf{w}) + \sum_{i=1}^r \lambda_i f_i(\mathbf{w}) + \sum_{i=1}^s \nu_i h_i(\mathbf{w}))$$

- $g$  is **concave** and can be  $-\infty$  for some  $\boldsymbol{\lambda}$  and  $\boldsymbol{\nu}$
- According to the constraints:  $f_i(x) \leq 0, h_i(x) = 0, \lambda_i > 0$ , we have:

$$f_0(\mathbf{w}^*) \geq L(\mathbf{w}^*, \boldsymbol{\lambda}, \boldsymbol{\nu}) \geq \inf_{\mathbf{w}} (L(\mathbf{w}^*, \boldsymbol{\lambda}, \boldsymbol{\nu})) = g(\boldsymbol{\lambda}, \boldsymbol{\nu})$$

- Definition of  $p^*$  (the value we want):

$$p^* = \inf\{f_0(\mathbf{w}) | f_i(\mathbf{w}) \leq 0, i = 1, \dots, r, h_i(\mathbf{w}) = 0, i = 1, \dots, s\}$$

- Then we have the following equation

$$d^* = \max_{\lambda \geq 0, \nu} g(\lambda, \nu) \leq \min f_0(\mathbf{w}^*) = p^*$$

- Lagrange dual problem (transformed by the Lagrange dual function)

$$\begin{aligned} & \max_{\lambda, \nu} g(\lambda, \nu) \\ & \text{subject to } \lambda \geq 0 \end{aligned}$$

- It is a convex optimization problem, with optimal value  $d^*$
- $d^*$  is the best lower bound on  $p^*$
- **Weak duality:**  $d^* \leq p^*$ , which always holds for convex and nonconvex problems
- **Strong duality:**  $d^* = p^*$  does not hold in general, but (usually) holds for convex problems, e.g., SVM

## Optimization of SVM: From Primal to Dual

- Primal Form:

- Notice that for constraints  $f_i$ , there is a set of constraints

$$\begin{aligned} & \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^m \xi_i \\ & \text{subject to } y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) \geq 1 - \xi_i, i = 1, \dots, m \\ & \quad \xi_i \geq 0, i = 1, \dots, m \end{aligned}$$

- Dual Form (Weak Duality):

$$\begin{aligned} L(\mathbf{w}, \mathbf{b}, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\beta}) &= \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i [y^{(i)} (\mathbf{w}^T \mathbf{x}^{(i)} + b) - 1 + \xi_i] - \sum_{i=1}^m \beta_i \xi_i \\ &= \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^m \alpha_i y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)} - \sum_{i=1}^m \alpha_i y^{(i)} b + \sum_{i=1}^m (C - \alpha_i - \beta_i) \xi_i + \sum_{i=1}^m \alpha_i \end{aligned}$$

- Notice that  $\alpha$  and  $\beta$  are two vectors contains the Lagrange Multiplier  $\lambda$ , there is no affine constraint
- We minimize the **Lagrangian** (definition of  $g$ ) w.r.t the primal variables  $\mathbf{w}, \mathbf{b}, \boldsymbol{\xi}$

$$\begin{aligned} \nabla_{\mathbf{w}} L &= \mathbf{w} - \sum_{i=1}^m \alpha_i y^{(i)} \mathbf{x}^{(i)} = 0 \quad (1) \\ \nabla_{\mathbf{b}} L &= - \sum_{i=1}^m \alpha_i y^{(i)} = 0 \quad (2) \\ \nabla_{\boldsymbol{\xi}} L &= C - \alpha_i - \beta_i = 0 \quad (3) \\ & \quad i = 1, \dots, m \end{aligned}$$

- Plugin the gradient equality to the Lagrangian to get the **Lagrange dual function**:

$$\begin{aligned}
g(\boldsymbol{\alpha}, \boldsymbol{\beta}) &= \inf_{\mathbf{w}, \mathbf{b}, \boldsymbol{\xi}} (L(\mathbf{w}, \mathbf{b}, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\beta})) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^m \alpha_i y^{(i)} \mathbf{w}^T \mathbf{x}^{(i)} \\
&\quad - \sum_{i=1}^m \alpha_i y^{(i)} b (= 0 \text{ according to (2)}) \\
&\quad + \sum_{i=1}^m (C - \alpha_i - \beta_i) \xi_i (= 0 \text{ according to (3)}) + \sum_{i=1}^m \alpha_i \\
&= \sum_{i=1}^m \alpha_i + \frac{1}{2} \left( \sum_{i=1}^m \alpha_i y^{(i)} \mathbf{x}^{(i)} \right)^T \left( \sum_{i=1}^m \alpha_i y^{(i)} \mathbf{x}^{(i)} \right) - \sum_{i=1}^m \alpha_i y^{(i)} \left( \sum_{i=1}^m \alpha_i y^{(i)} \mathbf{x}^{(i)} \right)^T \mathbf{x}^{(i)} \\
&\quad (\text{according to (1)}) \\
&= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j (\mathbf{x}^{(i)})^T \mathbf{x}^{(j)}
\end{aligned}$$

- The property:  $(a+b)^T(c+d) = (a^T+b^T)(c+d) = a^Tc + a^Td + b^Tc + b^Td$  is used in the last step
- Finally, we get the **dual form** of the primary form of SVM optimization:

$$\begin{aligned}
\max_{\boldsymbol{\alpha}} \quad & \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j (\mathbf{x}^{(i)})^T \mathbf{x}^{(j)} \\
\text{subject to } \nabla_{\mathbf{b}} L = & - \sum_{i=1}^m \alpha_i y^{(i)} = 0 \quad (2) \\
0 \leq \alpha_i \leq C, \quad & \alpha_i \geq 0, \beta_i \geq 0, \beta_i = C - \alpha_i \quad (3)
\end{aligned}$$

- Once the dual optimization is solved, we can find the  $\boldsymbol{\alpha}^*$ , the optimal  $\mathbf{w}^*$  can be computed by:

$$\mathbf{w}^* = \sum_{i=1}^m \alpha_i^* y^{(i)} \mathbf{x}^{(i)}$$

- How to compute the optimal bias  $b^*$

- For an item  $i$ :

$$\begin{aligned}
y^{(i)} ((\mathbf{w}^*)^T \mathbf{x}^{(i)} + b_i^*) &= 1 \\
(y^{(i)})^2 ((\mathbf{w}^*)^T \mathbf{x}^{(i)} + b_i^*) &= y^{(i)}, y \in \{-1, 1\} \\
(\mathbf{w}^*)^T \mathbf{x}^{(i)} + b_i^* &= y^{(i)} \\
b_i^* &= y^{(i)} - (\mathbf{w}^*)^T \mathbf{x}^{(i)}
\end{aligned}$$

- For the whole training set  $D$ :  $\sum_{i=1}^m b_i^* = \sum_{i=1}^m y^{(i)} - (\mathbf{w}^*)^T \mathbf{x}$
- Strong Duality Form:

- Complementary slackness: Assume the strong duality holds, and optimal variables are:  $\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*, \mathbf{w}^*, \mathbf{b}^*, \boldsymbol{\xi}^*$ , we have:

$$\begin{aligned}
p^* &= f_0(\mathbf{w}^*) = \inf_{\mathbf{w}} (f_0(\mathbf{w}) + \sum_{i=1}^r \lambda_i^* f_i(\mathbf{w}) + (\sum_{i=1}^s \nu_i^* h_i(\mathbf{w}) = 0)) \\
&\leq f_0(\mathbf{w}^*) + \sum_{i=1}^r \lambda_i^* f_i(\mathbf{w}^*) + (\sum_{i=1}^s \nu_i^* h_i(\mathbf{w}^*) = 0)
\end{aligned}$$

- $w^*$  minimize both  $f_0(\mathbf{w})$  and  $L(\mathbf{w}, \lambda^*, \nu^*)$  with Lagrange multiplier  $\lambda^*$  and  $\nu^*$
- $\sum_{i=1}^r \lambda_i^* f_i(\mathbf{w}^*) = 0$ , notice that  $\lambda_i^* f_i(\mathbf{w}^*) \leq 0$  for each  $i = 1, \dots, m$  (know as **complementary slackness**), which implies that  $\lambda_i^* f_i(\mathbf{w}^*) = 0$ . Therefore, there are only 2 situations:

- $\lambda_i > 0, f_i(\mathbf{w}^*) = 0; i = 1, \dots, m$
- $\lambda_i = 0, f_i(\mathbf{w}^*) < 0; i = 1, \dots, m$
- In the context of SVM with soft margin, we have the following complementary slackness (assume  $C > 0$ ):

$$\alpha_i^* [y^{(i)} ((x^{(i)})^T w^* + b^*) - 1 + \xi_i^*] = 0, \quad i = 1, \dots, m,$$

$$\beta_i^* \xi_i^* = 0, \quad i = 1, \dots, m.$$

Combining with

- $y^{(i)}((w^*)^T x^{(i)} + b^*) \geq 1 - \xi_i^*$ ,
- $\xi_i^* \geq 0$ ,
- $C - \alpha_i^* - \beta_i^* = 0$ ,
- $\alpha_i^* \geq 0$ ,
- $\beta_i^* \geq 0$ ,

we have

- if  $\alpha^* = 0$

$$\begin{cases} \alpha_i^* = 0 \\ C - \alpha_i^* - \beta_i^* = 0 \end{cases} \implies \beta_i^* = C, \quad (88)$$

$$\begin{cases} \beta_i^* = C \\ \beta_i^* \xi_i^* = 0 \end{cases} \implies \xi_i^* = 0, \quad (89)$$

$$\begin{cases} \xi_i^* = 0 \\ y^{(i)}((w^*)^T x^{(i)} + b^*) \geq 1 - \xi_i^* \end{cases} \implies y^{(i)}((w^*)^T x^{(i)} + b^*) \geq 1; \quad (90)$$

- if  $\alpha^* = C$

$$\begin{cases} \alpha_i^* = C \\ \alpha_i^* [y^{(i)} ((x^{(i)})^T w^* + b^*) - 1 + \xi_i^*] = 0 \end{cases} \implies y^{(i)} ((x^{(i)})^T w^* + b^*) - 1 + \xi_i^* = 0, \quad (91)$$

$$\begin{cases} y^{(i)} ((x^{(i)})^T w^* + b^*) - 1 + \xi_i^* = 0 \\ \xi_i^* \geq 0 \end{cases} \implies y^{(i)} ((x^{(i)})^T w^* + b^*) \leq 1; \quad (92)$$

- if  $0 < \alpha^* < C$

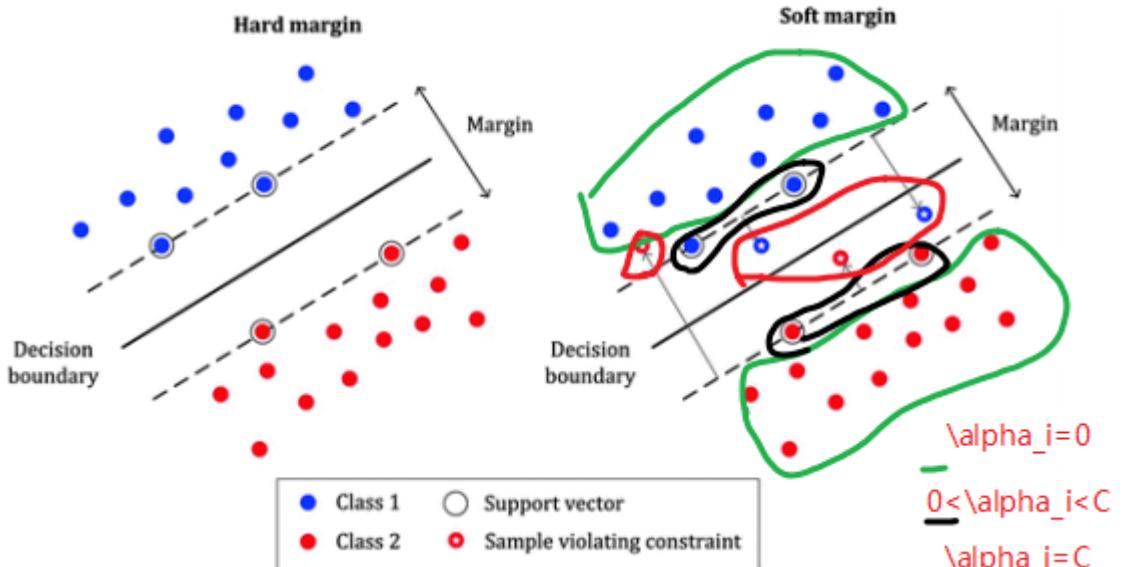
$$\begin{cases} 0 < \alpha^* < C \\ C - \alpha_i^* - \beta_i^* = 0 \end{cases} \implies \beta_i^* > 0, \quad (93)$$

$$\begin{cases} 0 < \alpha^* < C \\ \alpha_i^* [y^{(i)}((x^{(i)})^T w^* + b^*) - 1 + \xi_i^*] = 0 \end{cases} \implies y^{(i)}((x^{(i)})^T w^* + b^*) - 1 + \xi_i^* = 0, \quad (94)$$

$$\begin{cases} \beta_i^* > 0 \\ \beta_i^* \xi_i^* = 0 \end{cases} \implies \xi_i^* = 0, \quad (95)$$

$$\begin{cases} \xi_i^* = 0 \\ y^{(i)}((w^*)^T x^{(i)} + b^*) - 1 + \xi_i^* = 0 \end{cases} \implies y^{(i)}((w^*)^T x^{(i)} + b^*) = 1. \quad (96)$$

- Support Vectors:



- Why Bother SVM Dual

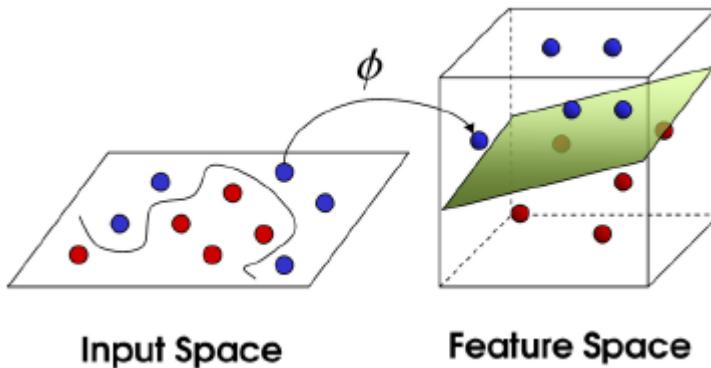
- How would one classify a new point  $x$  with primal SVM solved
  - We need to explicitly compute the scalar product  $\mathbf{w}^T x + b$
- How would one classify a new point  $x$  with dual SVM solved
  - Very efficient as we only compute a few “support vectors” whose  $\alpha_i$ ’s are nonzero (for zero  $\alpha_i$ , the constraints are not used)
  - In this situation, weight matrix can be considered as a linear combination of training samples

$$\begin{aligned}
\mathbf{w}^T \mathbf{x} + b &= \left( \sum_{i=1}^m \alpha_i y^{(i)} \mathbf{x}^{(i)} \right)^T \mathbf{x} + b \\
&= \sum_{i=1}^m \alpha_i y^{(i)} (\mathbf{x}^{(i)})^T \mathbf{x} + b
\end{aligned}$$

- More importantly, the dual form lends itself easily to the *kernel trick*

## Kernel Methods

- The drawback of Linear Models
  - Decision boundaries are by definition linear in the input feature space
  - High bias on complex data
  - Need to **relax the constraint of linear decision boundaries** while **retaining the nice properties of linear classifiers** => increase the dimension of the feature space



- Basis Function Expansion
    - Simple method to increase dimension: Apply a set (with cardinality  $l$ ) of function  $\{\phi_1, \dots, \phi_l\}$  to the raw feature vector  $\mathbf{x}$  to map it in to a new feature space:
- $$\phi(\mathbf{x}) = [\phi_1(\mathbf{x}), \dots, \phi_l(\mathbf{x})]^T$$
- This is called a *basis function expansion* since  $l > n$  in general.
  - This requires that we know the functions  $\phi_1, \dots, \phi_l$  in advance
  - We then define a linear classifier (SVM or logistic regression) in this new feature space:

$$\mathbf{w}^T \phi(\mathbf{x}) + b$$

- Basis Function Expansion Examples
  - Degree 2 Polynomial Basis:** We include all single features  $x_j$ , their squares  $x_j^2$ , and all products of two distinct features  $x_i x_j$
  - Degree  $d$  Polynomial Basis:** We include all single features  $x_j$ , and all unique products of between 2 and  $d$  features
  - The problem is that the space complexity of representing the expanded set of features is essentially  $O(n^d)$  => Need to be simplify
- Representer Theorem
  - One of the interesting properties of SVMs is that the optimal weight vectors can always be expressed as a weighted linear combination of the data vectors:

$$\mathbf{w} = \sum_{i=1}^m \gamma_i x^{(i)}$$

- This result is called the *representer theorem*
- Dependence on Inner products
  - Plugging this result back in to the SVM objective, we find that the objective only depends on the data through inner products
  - For normalized input vector  $\mathbf{x}, \mathbf{z}$ ,  $\mathbf{x}^T \mathbf{z} = \|\mathbf{x}\| \cdot \|\mathbf{z}\| \cos\theta = \cos\theta$ , which represents the similarity between the two vectors

$$\begin{aligned} & \min_{\gamma, b, \xi} \frac{1}{2} \sum_{j=1}^m \sum_{k=1}^m \gamma_j \gamma_k \phi(\mathbf{x}^{(j)})^T \phi(\mathbf{x}^{(k)}) + C \sum_{i=1}^m \xi_i \\ & \text{subject to } y^{(i)} \left( \left( \sum_{j=1}^m \gamma_j \phi(\mathbf{x}^{(j)}) \right)^T \phi(\mathbf{x}^{(i)}) + b \right) \geq 1 - \xi_i, \\ & \quad \xi_i \geq 0, i = 1, \dots, m \end{aligned}$$

- The Kernel Trick
  - Amazingly, for many useful basis function expansions  $\phi(x)$ , it is possible to find a function  $K(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x})^T \phi(\mathbf{z})$  that can **compute the inner product under the basis expansion without ever explicitly performing the basis function expansion**, where  $\mathbf{x}, \mathbf{z} \in \mathbb{R}^n$  are two input feature vectors
  - Such functions are called *kernel functions* and this is known as the “Kernel Trick”
  - In other words, we can also directly learn the parameters  $\gamma_i$ ’s using the kernel trick, without constructing the basis expansion (since the expression  $\mathbf{w}^T \phi(\mathbf{x}) + b$  also needs inner product)
  - Interestingly, there **exist kernels for which the basis function expansion implied by the kernel is infinite dimensional**
- Examples of Kernel Functions
  - **Degree  $d$  Polynomial Kernel:**  $K_p(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z} + c)^d$ , which corresponds to a feature mapping to an  $\binom{d+n}{d}$  feature space. Where  $\mathbf{x}, \mathbf{z} \in \mathbb{R}^n$  are two input feature vectors, and  $c \geq 0$  is a constant. If  $c = 0$ ,  $K(\mathbf{x}, \mathbf{z})$  reduces to the homogeneous(同質) polynomial kernel
    - proof:
      - Let  $k_0, k_1, \dots, k_n$  denote non-negative integers, such that  $\sum_{j=0}^n k_j = d$ . The multinomial expansion(多項式展開) of the inhomogeneous kernel is then given as

$$\begin{aligned} K(\mathbf{x}, \mathbf{z}) &= (\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z} + c)^d = (c + \sum_{j=1}^n x_j z_j)^d = (c + x_1 z_1 + \dots + x_n z_n)^d \\ &= \sum_{k_0+\dots+k_n=d} \binom{d}{k_0, \dots, k_n} c^{k_0} (x_1 z_1)^{k_1} \dots (x_n z_n)^{k_n} \\ &, (\text{where } \binom{d}{k_0, \dots, k_n} \text{ means the joint selection number}) \\ &= \sum_{k_0+\dots+k_n=d} \binom{d}{k_0, \dots, k_n} c^{k_0} (x_1^{k_1} \dots x_n^{k_n})(z_1^{k_1} \dots z_n^{k_n}) \\ &= \sum_{k_0+\dots+k_n=d} \left( \sqrt{\binom{d}{k_0, \dots, k_n}} c^{k_0} \prod_{j=1}^n x_j^{k_j} \right) \left( \sqrt{\binom{d}{k_0, \dots, k_n}} c^{k_0} \prod_{j=1}^n z_j^{k_j} \right) \end{aligned}$$

- Then, the mapping  $\phi : \mathbb{R}^n \mapsto \mathbb{R}^l$  is given as the vector:

$$\phi(\mathbf{x}) = [\dots, \sqrt{\binom{d}{k_0, \dots, k_n}} c^{k_0} \prod_{j=1}^n x_j^{k_j}, \dots]^T,$$

- where  $(k_0, \dots, k_n)$  ranges over all the possible assignments, such that  $k_j \geq 0$  for  $0 \leq j \leq n$  and  $\sum_{j=0}^n k_j = d$
- i.e. partition  $d$  items to  $n+1$  (one bias) buckets => select the separator =>  $\binom{d+n}{n} = \binom{d+n}{d}$
- **Gaussian/RBF Kernel:**  $K_G(\mathbf{x}, \mathbf{z}) = \exp(-\gamma \|\mathbf{x} - \mathbf{z}\|_2^2)$ , which corresponds to an infinite dimensional feature space

- proof:

- First, use Taylor expansion of  $f(x) = e^x$  at the point  $(0, 1)$ :

$$\begin{aligned} f(\theta) = e^\theta &= \sum_{d=0}^{\infty} \frac{(\theta - 0)^d}{d!} f^d(0) = \sum_{d=0}^{\infty} \frac{\theta^d}{d!} \\ &= 1 + \theta + \frac{1}{2!}\theta^2 + \frac{1}{3!}\theta^3 + \dots \end{aligned}$$

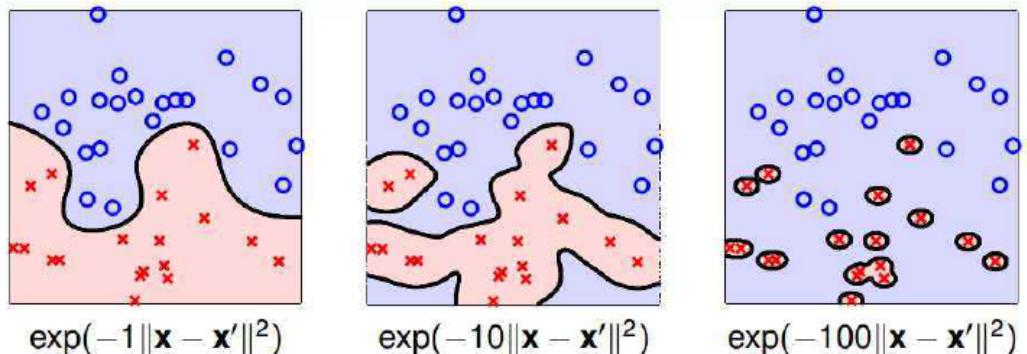
- Further, noting that  $\|\mathbf{x} - \mathbf{z}\|_2^2 = \|\mathbf{x}\|_2^2 + \|\mathbf{z}\|_2^2 - 2\mathbf{x}^T \mathbf{z}$ , we can rewrite the Gaussian kernel as follows:

$$\begin{aligned} K(\mathbf{x}, \mathbf{z}) &= \exp\{-\gamma \|\mathbf{x} - \mathbf{z}\|_2^2\} \\ &= \exp\{-\gamma \|\mathbf{x}\|_2^2\} \exp\{-\gamma \|\mathbf{z}\|_2^2\} \exp\{2\gamma \mathbf{x}^T \mathbf{z}\} \\ &= \exp\{-\gamma \|\mathbf{x}\|_2^2\} \exp\{-\gamma \|\mathbf{z}\|_2^2\} \sum_{d=0}^{\infty} \frac{(2\gamma)^d}{d!} (\mathbf{x}^T \mathbf{z})^d \\ &= \exp\{-\gamma \|\mathbf{x}\|_2^2\} \exp\{-\gamma \|\mathbf{z}\|_2^2\} \sum_{d=0}^{\infty} \frac{(2\gamma)^d}{d!} \left( \sum_{k_1+\dots+k_n=d} \binom{d}{k_0, \dots, k_n} c^{k_0} \prod_{j=1}^n (x_j z_j)^{k_j} \right) \\ &= \sum_{d=0}^{\infty} \sum_{k_1+\dots+k_n=d} (\omega \times \exp\{-\gamma \|\mathbf{x}\|_2^2\} \prod_{j=1}^n (x_j)^{k_j}) (\omega \times \exp\{-\gamma \|\mathbf{z}\|_2^2\} \prod_{j=1}^n (z_j)^{k_j}) \\ &\quad \text{where } \omega = \sqrt{\left( \sum_{k_0, \dots, k_n} \binom{d}{k_0, \dots, k_n} \frac{(2\gamma)^d}{d!} \right)} \end{aligned}$$

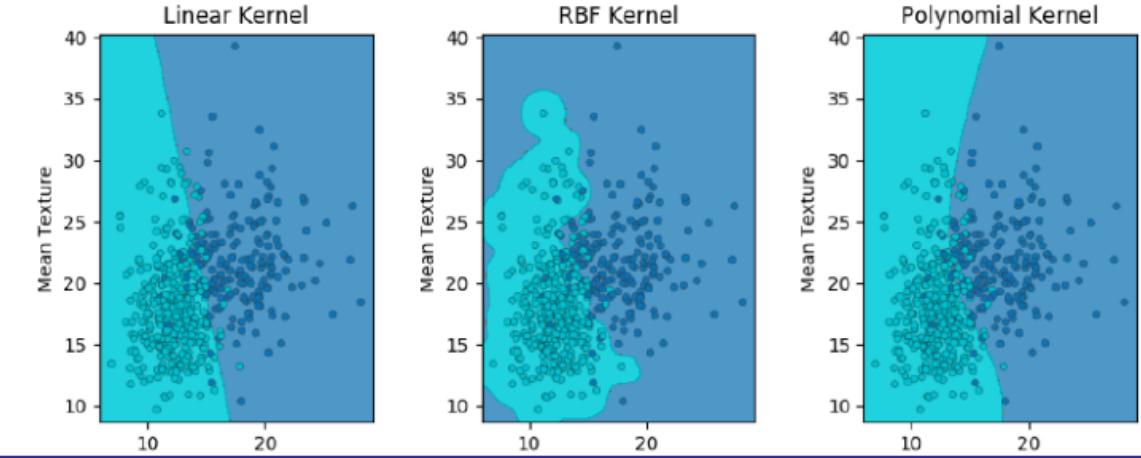
- Where the mapping  $\phi : \mathbb{R}^n \mapsto \mathbb{R}^\infty$  is given as the vector

$$\phi(\mathbf{x}) = [\dots, \sqrt{\left( \sum_{k_0, \dots, k_n} \binom{d}{k_0, \dots, k_n} \frac{(2\gamma)^d}{d!} \right)} \exp\{-\gamma \|\mathbf{x}\|_2^2\} \prod_{j=1}^n (x_j)^{k_j}, \dots]^T$$

- Since  $d$  goes to  $\infty$  and each  $d$  have  $\binom{d+n}{d}$  combinations, the  $\phi$  maps the input space into an infinite dimensional feature space, we obviously cannot compute  $\phi(\mathbf{x})$ , yet computing the Gaussian kernel  $K(\mathbf{x}, \mathbf{z})$  is straightforward and efficient
- $\gamma$  also matters since it controls how the higher level Taylor expansions expressed



- The graphs:



- Kernel Trick for other ML methods

- Many learning methods depend on computing inner products between features, e.g., linear regression
- For those methods, we can use a kernel function in the place of the inner products, i.e., "kernelizing" the methods, thus, introducing nonlinear features/basis
- Instead of first transforming the original features into the new feature space and then computing the inner product, we can compute the inner product in the new feature space directly through the kernel function

- Kernel Ridge Regression

- Recall ridge regression:  $\min_w \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2$

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_{n+1})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{X}^T (\mathbf{X} \mathbf{X}^T + \lambda \mathbf{I}_m)^{-1} \mathbf{y}$$

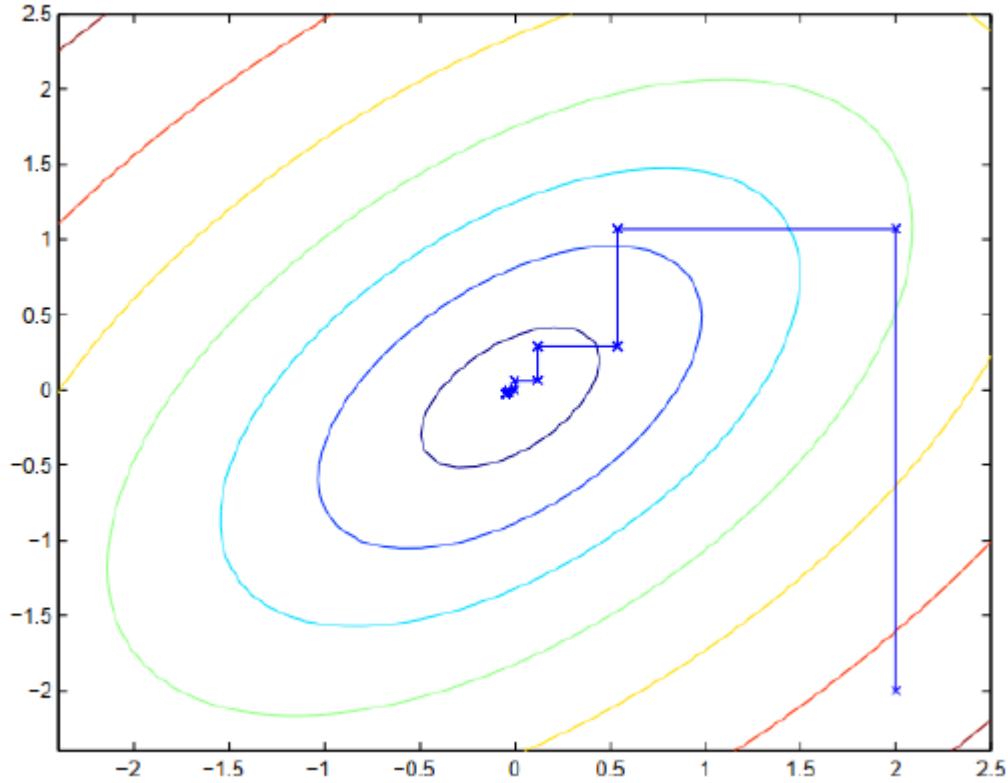
$$\begin{aligned}
& (\mathbf{P} + \mathbf{B}^T \mathbf{R}^{-1} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{R}^{-1} \\
&= (\mathbf{P} + \mathbf{B}^{-1} \mathbf{R}(\mathbf{B}^T)^{-1}) \mathbf{B}^T \mathbf{R}^{-1} \\
&\stackrel{(P^{-1} = \lambda I_{n+1}, B = X, R = I_m)}{=} \mathbf{P} \mathbf{B}^T \mathbf{B}^{-1} + \mathbf{P} \mathbf{B}^T (\mathbf{B}^T)^{-1} \mathbf{P}^T \mathbf{B}^{-1} \\
&= \mathbf{P} \mathbf{B}^T (\mathbf{R}^{-1} + (\mathbf{B}^T)^T \mathbf{P}^{-1} \mathbf{B}^{-1}) \\
&= \mathbf{P} \mathbf{B}^T (\mathbf{R} + \mathbf{B} \mathbf{P} \mathbf{B}^T)^{-1} \\
&\quad \text{For } \underbrace{(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_{n+1})^{-1}}_{U} \mathbf{X}^T \mathbf{y} \\
&\quad \stackrel{(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_{n+1})^{-1} = \frac{1}{\lambda} \mathbf{X}^T (\mathbf{I}_m + \frac{1}{\lambda} \mathbf{X} \mathbf{X}^T)^{-1}}{=} \frac{1}{\lambda} \mathbf{X}^T (\mathbf{I}_m + \frac{1}{\lambda} \mathbf{X} \mathbf{X}^T)^{-1} \\
&\quad \stackrel{(\mathbf{I}_m + \frac{1}{\lambda} \mathbf{X} \mathbf{X}^T)^{-1} = \frac{\lambda}{\lambda + \mathbf{X}^T \mathbf{X}}} {=} \mathbf{X}^T (\lambda \mathbf{I}_m + \mathbf{X} \mathbf{X}^T)^{-1} \\
&\quad \stackrel{(\lambda \mathbf{I}_m + \mathbf{X} \mathbf{X}^T)^{-1} = \mathbf{X}^T (\lambda \mathbf{I}_m + \mathbf{X} \mathbf{X}^T)^{-1}}{=} \mathbf{X}^T (\lambda \mathbf{I}_m + \mathbf{X} \mathbf{X}^T)^{-1} \mathbf{y} \\
&= \mathbf{X}^T (\lambda \mathbf{I}_m + \mathbf{X} \mathbf{X}^T)^{-1} \mathbf{y}.
\end{aligned}$$

- This equation can be rewritten as:  $\mathbf{w}^* = \sum_{i=1}^m \alpha_i \mathbf{x}^{(i)}$  with  $\boldsymbol{\alpha} = (\mathbf{X} \mathbf{X}^T + \lambda \mathbf{I}_m)^{-1} \mathbf{y}$ . This indicates that the solution  $\mathbf{w}^*$  must lie in the span of the data cases, which makes intuitive sense because the algorithm is linear in data space
- For any new point  $\mathbf{x}$ , we make prediction by:  $y = \mathbf{w}^T \mathbf{x} = \mathbf{y}^T (\mathbf{X} \mathbf{X}^T + \lambda \mathbf{I}_m)^{-1} \mathbf{X} \mathbf{x}$ , where  $(\mathbf{X} \mathbf{X}^T)^T = \mathbf{X} \mathbf{X}^T$
- Kernelizing the method, we have  $y = \mathbf{w}^T \mathbf{x} = \mathbf{y}^T (\mathbf{K} + \lambda \mathbf{I}_m)^{-1} \kappa(\mathbf{x})$ , where  $K_{ij} = \mathbf{K}(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \phi(\mathbf{x}^{(i)})^T \phi(\mathbf{x}^{(j)})$ ,  $\kappa_i(\mathbf{x}) = \mathbf{K}(\mathbf{x}^{(i)}, \mathbf{x})$  and  $\phi(\mathbf{x}) : \mathbb{R}^n \mapsto \mathbb{R}^l$  is a feature mapping

- No extra computational burden is incurred, but it now becomes a (much more powerful) nonlinear regression model
- Trade-Offs
  - Linear SVM and logistic regression are both discriminative linear classifiers with identical capacity and space complexity. They have very similar convex loss functions and identical regularizers. They are fast and space efficient, but can have high bias
  - Basis expansion requires more space ( $O(ml)$  ( $l$  functions) for data and  $O(l)$  for parameters), but yields non linear classifiers that have lower bias
  - Kernel SVM requires  $O(m^2)$ (a function with two parameters, takes all the combination among two data vectors) space for storing all the kernel values during training and have  $O(m)$  parameters. This can still be much lower than  $O(ml)$  for large sets of basis functions. Kernels also yield non-linear classifiers that have lower bias (How does the complexity calculated)
  - Kernel SVM often has at least two **hyperparameters** (C and a kernel hyperparameter). These need to be set jointly, which may be computationally expensive
  - Gaussian kernel SVM has infinite capacity, closely related to KNN
  - **Unlike logistic loss, hinge loss is not differentiable, so SVMs require more advanced optimization methods**
  - It is somewhat more difficult to form a multi-class SVM than a multi-class logistic regression model, and many implementations use “hacks” like one-vs-all

## Sequential Minimal Optimization

- Solving SVM: Coordinate Descent
  - Another method of optimization (compared with GD)
  - Consider the unconstrained optimization problem:  $\min(\theta_1, \theta_2, \dots, \theta_n)$ 
    - The coordinate descent method first selects an initial point  $\theta^{(0)} \in \mathbb{R}^n$  and repeats:
      - Choose an index  $i$  from 1 to  $n$
      - $\theta_i^{(t+1)} = \operatorname{argmin}_{\hat{\theta}_i} l(\theta_1^{(t)}, \dots, \theta_{i-1}^{(t)}, \hat{\theta}_i, \theta_{i+1}^{(t)}, \dots, \theta_n^{(t)})$
    - At each iteration, the algorithm determines a coordinate, and minimizes over the corresponding direction while fixing all other coordinates
    - Coordinate descent is applicable in both differentiable and derivative-free contexts



- Solving SVM: Sequential Minimal Optimization (SMO)
  - SMO is an algorithm for solving the quadratic programming problem that arises during the training of SVMs, invented by John C. Platt in 1998
  - Requirements:
    - Sequential:
      - Not parallel
      - Optimize a set of two Lagrange multiplier
    - Minimal : Optimize the smallest possible sub-problem at each step
    - Optimization: Satisfy the constraint for the **chosen pair** of Lagrange multipliers
  - Recall the dual form of kernelized SVM with soft margin:

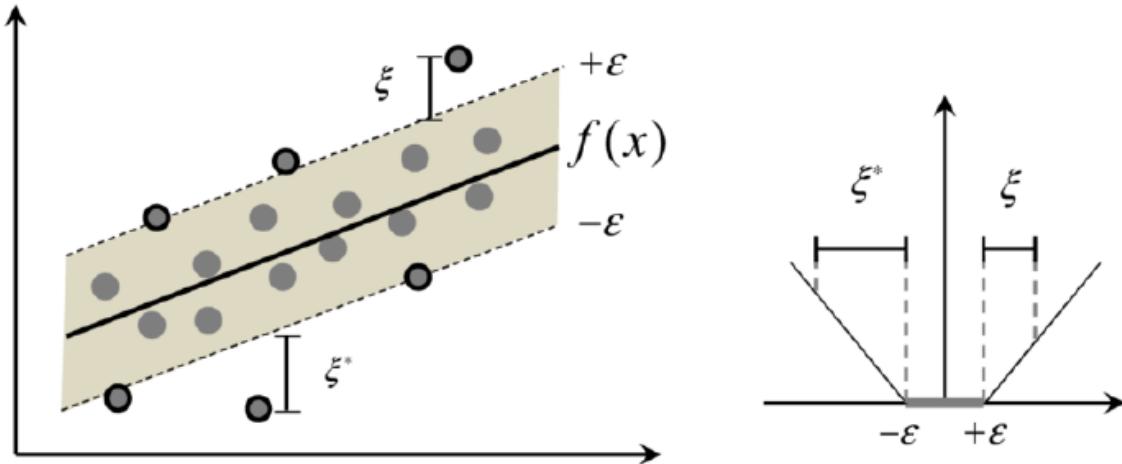
$$\begin{aligned}
 & \max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) \\
 & \text{s.t. } \sum_{i=1}^m \alpha_i y^{(i)} = 0, \\
 & \quad 0 \leq \alpha_i \leq C, i = 1, \dots, m
 \end{aligned}$$

- SMO is derived by decomposing the problem to its extreme and optimizing **a minimal subset of just two data points** at each iteration
  - The sub-problem for two data points admits an analytical solution, eliminating the need to use an iterative quadratic programming optimizer as part of the algorithm
  - To satisfy the second constraint, we need to use two data points, if use one data point, the second constraint will be broken since one of  $\alpha_i$  will be changed while others are not => when one  $\alpha$  is changed, the other one will be changed to fit the constraint

- At each step SMO chooses two elements  $\alpha_i$  and  $\alpha_j$  to jointly optimize, finds the optimal values for those two parameters given that all the others are fixed, and updates the  $\alpha$  vector accordingly
- The choice of the two points is determined by a **heuristic**, while the optimization of the two multipliers is performed analytically
- Despite needing more iterations to converge, each iteration uses so few operations that the algorithm exhibits an overall speed-up of some orders of magnitude

## Support Vector Regression (SVR)

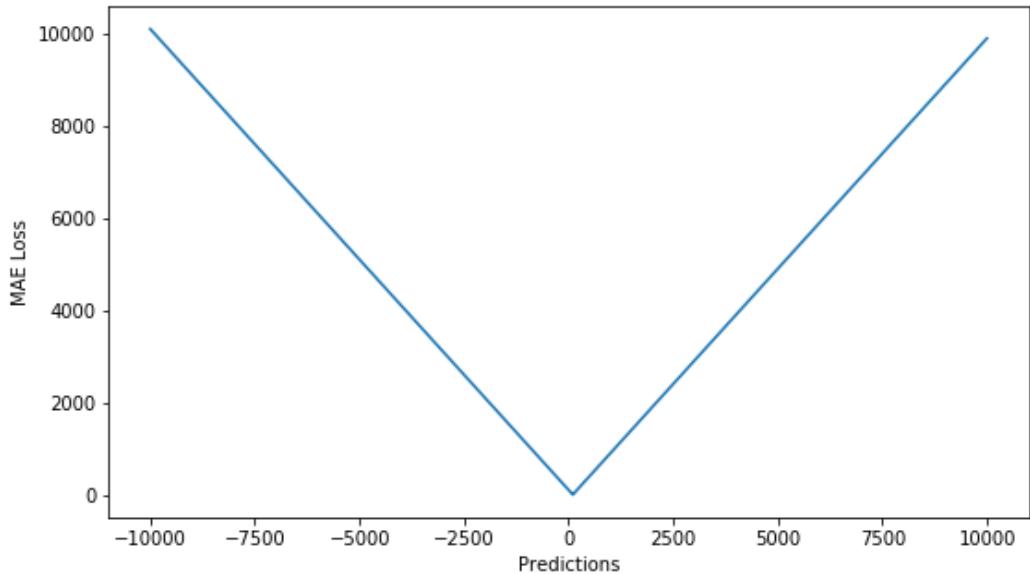
- Support vector regression (SVR) is the generalization of SVM to the case of regression
- Use a different objective function. In the following case, the *epsilon insensitive loss*



- Where MAE is like this:

$$MAE = \frac{\sum_{i=1}^n |y_i - y_i^p|}{n}$$

Range of predicted values: (-10,000 to 10,000) | True value: 100



Plot of MAE Loss (Y-axis) vs. Predictions (X-axis)

- Primal Form:

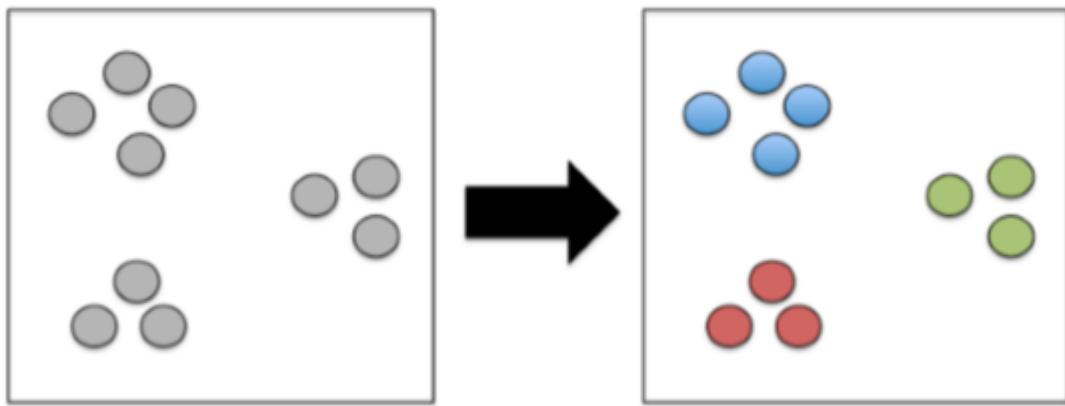
$$\begin{aligned} & \min_{\mathbf{w}, b, \xi, \xi^*} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*) \\ & \text{subject to } y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)} - b \leq \epsilon + \xi_i, i = 1, \dots, m \\ & \quad \mathbf{w}^T \mathbf{x}^{(i)} + b - y^{(i)} \leq \epsilon + \xi_i^*, i = 1, \dots, m \\ & \quad \xi_i, \xi_i^* \geq 0, i = 1, \dots, m \end{aligned}$$

- Where  $\xi_i$  and  $\xi_i^*$  represent different level of tolerance of different side of the decision boundary for the data point  $i$
- It can be generalized to non-linear models using the same representer theorem as in SVM (i.e. the kernel trick)
- It is a quadratic programming problem that can be efficiently solved by SMO
- Trade-Offs
  - SVR is more robust to outliers than OLS (which minimizing the MSE) due to the loss being linear in the tails instead of quadratic. This is related to a long line of work on robust regression
  - Kernel SVR has low bias and good capacity control, **but use of cross validation to select regularization hyperparameters is critical**
  - The kernel matrix computation is quadratic in the data dimension, and the model has a support vector property

# Lecture 08 K-Means Clustering and Gaussian Mixture Models

## Introduction

- The Classification Task: Given a feature vector  $\mathbf{x} \in \mathbb{R}^n$  that describes an object that belongs to one of  $|Y|$  classes from the set  $Y$ , predict which class the object belongs to
- The Clustering Task: Give a collection of data cases  $\mathbf{x}^{(i)} \in \mathbb{R}^n$ , partition the data cases into groups such that the data cases within each partition are more similar to each other than they are to data cases in other partitions



## Exhaustive Clustering

- Define a Clustering
  - Suppose we have  $m$  data cases  $D = \{\mathbf{x}^{(i)}\}_{i=1,\dots,m}$
  - A clustering of the  $m$  cases into  $K$  clusters is a partitioning of  $D$  into  $K$  mutually disjoint subsets  $C = C_1, \dots, C_k$  such that  $C_1 \cup \dots \cup C_K = D$
  - Suppose we have a function  $f(C)$  that takes a partitioning  $C$  of the data set  $D$  and returns a score with lower scores indicating better clustering
  - The optimal clustering according to  $f$  is simply given by:  $\text{argmin}_C f(C)$
  - The complexity of exhaustive clustering depends on the total number of partitions
- Number of Clusterings
  - The total number of clusterings of a data set with  $m$  elements is the Bell number  $B_m$ , where  $B_0 = B_1 = 1$  and  $B_{m+1} = \sum_{k=0}^m \binom{m}{k} B_k$
  - Wrong: The recursive relation means that, for the first cluster we choose  $m - k$  elements, and calculate how many ways of combination for  $k$  elements left, and sum all the possibility of  $k$ 's (It is wrong because for a partition, the order of choice doesn't matter)
  - True proof: Let  $x^{(1)}$  to be a fixed to the first cluster, we can choose  $m - k$  elements in the same cluster with  $x^{(1)}$ , then the order is fixed. Therefore, the equation is for  $m + 1$  instead of  $m$
  - The complexity of exhaustive clustering scales with  $B_m$  and is thus computationally totally intractable for general scoring functions

- We will need either approximation algorithms or scoring functions with special properties

## K-Means Clustering

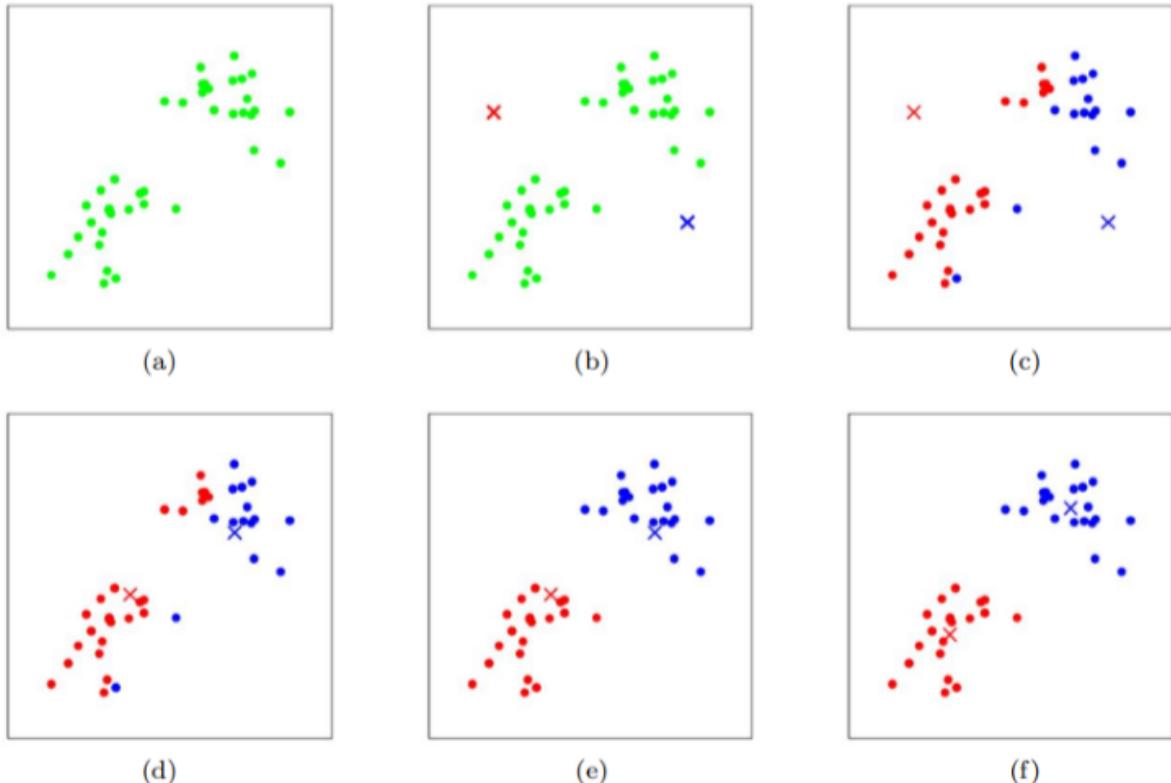
- The K-means algorithm is an iterative optimization algorithm for clustering that alternates between two steps
- The algorithm maintains a set of  $K$  cluster centroids or prototypes  $\{\mu_k\}$  that represent the average (mean) feature vector of the data cases in each cluster
- Algorithm: Suppose we let  $z^{(i)}$  indicate which cluster  $\mathbf{x}^{(i)}$  belongs to and  $\mu_k \in \mathbb{R}^n$  be the cluster centroid/prototype for cluster  $k$ . The two main steps of the algorithm can then be expressed as follows:
  - In the first step, the distance between each data case and each prototype is computed, and each data case is assigned to the nearest prototype

$$1 \quad z^{(i)} = \arg \min_k \|\mu_k - \mathbf{x}^{(i)}\|_2^2 \text{ (the assignment step)}$$

- In the second step, the prototypes are updated to the mean of the data cases assigned to them

$$2 \quad \mu_k = \frac{\sum_{i=1}^m \mathbb{I}[z^{(i)} = k] \mathbf{x}^{(i)}}{\sum_{i=1}^m \mathbb{I}[z^{(i)} = k]} \text{ (the update step)}$$

- Example:



- The K-Means Objective
  - The  $K$ -means algorithm attempts to minimize the sum of the within-cluster variation over all clusters (also called the within-cluster sum of squares):

$$\min_{\mathbf{z}, \boldsymbol{\mu}} \ell(\mathbf{z}, \boldsymbol{\mu}) = \min_{\mathbf{z}, \boldsymbol{\mu}} \sum_{i=1}^m \|\mathbf{x}^{(i)} - \boldsymbol{\mu}_{z^{(i)}}\|^2$$

- It can be shown that  $K$ -means is exactly coordinate descent on  $\ell$ . Specifically, the assignment step minimizes  $\ell$  with respect to  $z$  while holding  $\boldsymbol{\mu}$  fixed, and the update step minimizes  $\ell$  with respect to  $\boldsymbol{\mu}$  while holding  $z$  fixed. Thus,  $\ell$  must monotonically decrease, and the value of  $\ell$  must converge
- Note that  $\ell$  has many local optima in general, each corresponding to a different clustering of the data. Finding the global optimum is not computationally tractable => **highly sensitive to initialization**
- Initialization
  - It is common to perform multiple random re-starts of the algorithm and take the clustering with the minimal total variation
  - Common initializations include setting the initial centers to be randomly selected data cases, setting the initial partition to a random partition, and selecting centers using a “furthest first”-style heuristic (more formally known as K-means++)
  - It often helps to initially run with  $K \log(K)$  clusters, then merge clusters to get down to  $K$  and run the algorithm from that initialization
- Issues
  - Only works with Euclidean distance. An alternative based on Manhattan distance is called the K-medians algorithm
  - Pre-processing like re-scaling/normalizing features can completely change the results
  - We need some way to determine the “right” number of clusters to focus on. We want to cluster on salient differences between data cases, not noise
  - The number of iterations to convergence is often small (like 20), but examples can be constructed that require an exponential number of steps to converge
  - Results in a hard assignment of data cases to clusters, which may be a problem if there are outliers

## Expectation Maximization

- EM solves a maximum likelihood problem of the form:

$$L(\boldsymbol{\theta}) = \sum_{i=1}^m \log p(\mathbf{x}^{(i)}; \boldsymbol{\theta}) = \sum_{i=1}^m \log \sum_{z^{(i)} \in Z} p(\mathbf{x}^{(i)}, z^{(i)}; \boldsymbol{\theta})$$

- $\boldsymbol{\theta}$  : Parameters of the probabilistic model we try to find
- $\{\mathbf{x}^{(i)}\}_{i=1,\dots,m}$  : Observed training examples
- $\{z^{(i)}\}_{i=1,\dots,m}$  : Unobserved latent variables
- $Z$ : Finite set of categorical values that  $z^{(i)}$  can take on
- EM Derivation

$$\begin{aligned}
& \sum_{i=1}^m \log \sum_{z^{(i)} \in Z} p(\mathbf{x}^{(i)}, z^{(i)}; \boldsymbol{\theta}) \\
&= \sum_{i=1}^m \log \sum_{z^{(i)} \in Z} q(z^{(i)}) \frac{p(\mathbf{x}^{(i)}, z^{(i)}; \boldsymbol{\theta})}{q(z^{(i)})}, \text{ (where } q \text{ is an arbitrary distribution on } z \text{ )} \\
&= \sum_{i=1}^m \log \mathbb{E}_{z^{(i)} \sim q} \left[ \frac{p(\mathbf{x}^{(i)}, z^{(i)}; \boldsymbol{\theta})}{q(z^{(i)})} \right] \geq \sum_{i=1}^m E_{z^{(i)} \sim q} [\log \frac{p(\mathbf{x}^{(i)}, z^{(i)}; \boldsymbol{\theta})}{q(z^{(i)})}] \\
&= \sum_{i=1}^m \sum_{z^{(i)} \in Z} q(z^{(i)}) \log p(\mathbf{x}^{(i)}, z^{(i)}; \boldsymbol{\theta}) - \sum_{i=1}^m \sum_{z^{(i)} \in Z} q(z^{(i)}) \log q(z^{(i)})
\end{aligned}$$

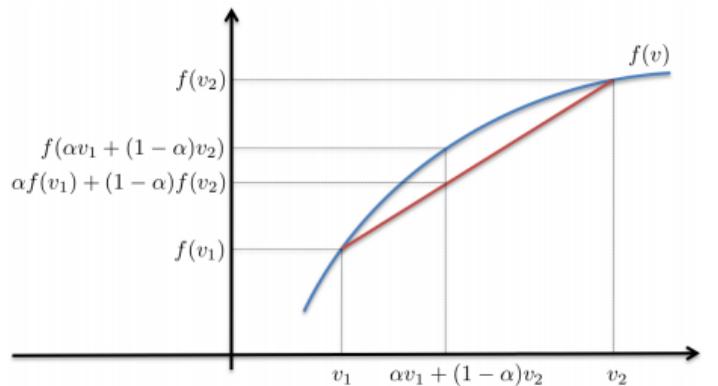
- Jensen's inequality:

- Suppose  $f : \mathbb{R} \mapsto \mathbb{R}$  is concave, then for all probability distributions  $p$  and all probability distributions  $q$  and all functions  $g : \mathbb{R} \mapsto \mathbb{R}$ , we have:

$$f(\mathbb{E}_{v \sim p}[g(v)]) \geq E_{v \sim p}[f(g(v))]$$

with equality iff  **$f$  is an affine function or  $g(v)$  is a constant**

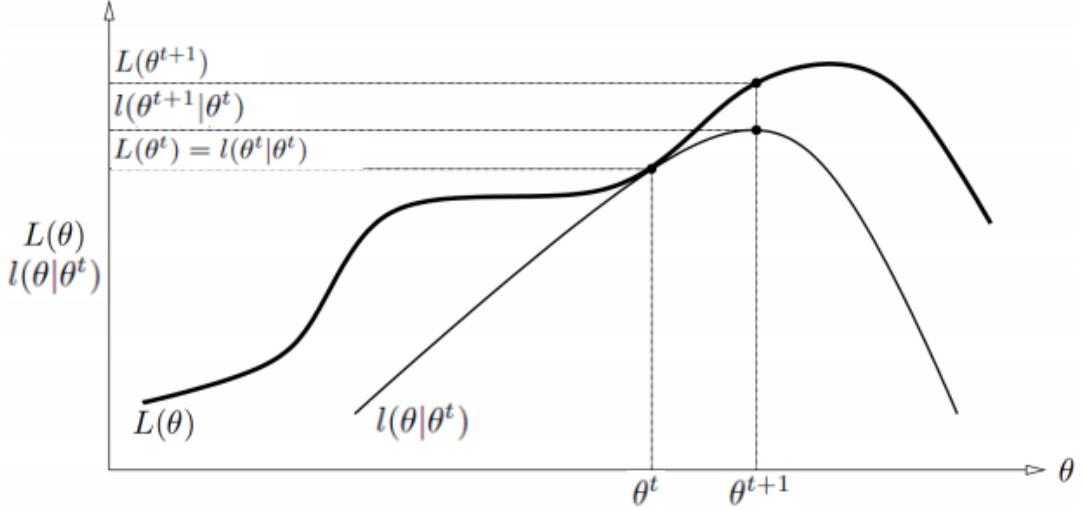
**Illustration:**  
 $p(V = v_1) = \alpha$ ,  
 $p(V = v_2) = 1 - \alpha$



- Lower bound:

$$\begin{aligned}
L(\boldsymbol{\theta}) &= \sum_{i=1}^m \log \sum_{z^{(i)} \in Z} p(\mathbf{x}^{(i)}, z^{(i)}; \boldsymbol{\theta}) \\
&\geq \sum_{i=1}^m \sum_{z^{(i)} \in Z} q(z^{(i)}) \log p(\mathbf{x}^{(i)}, z^{(i)}; \boldsymbol{\theta}) - \sum_{i=1}^m \sum_{z^{(i)} \in Z} q(z^{(i)}) \log q(z^{(i)}) \\
&= l(\boldsymbol{\theta})
\end{aligned}$$

- Since  $f(x) = \log(x)$ ,  $f$  is impossible to be an affine, the equality holds iff  $g(q) = \frac{p(\mathbf{x}^{(i)}, z^{(i)}; \boldsymbol{\theta})}{q(z^{(i)})}$  is a constant, which can be achieved when  $q(z^{(i)}) = p(z^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta}) \propto p(\mathbf{x}^{(i)}, z^{(i)}; \boldsymbol{\theta})$
- The EM algorithm repeatedly carries out the following two steps until convergence:
  - E-step: For each  $i$ , compute  $q(z^{(i)}) = p(z^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta})$ 
    - We do not fill in the unobserved  $z^{(i)}$  with hard values, but find a posterior distribution  $q(z^{(i)})$ , given  $\mathbf{x}^{(i)}$  and  $\boldsymbol{\theta}$
  - M-step:  $\boldsymbol{\theta} = \arg \max_{\boldsymbol{\theta}} l(\boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^m \sum_{z^{(i)} \in Z} q(z^{(i)}) \log p(\mathbf{x}^{(i)}, z^{(i)}; \boldsymbol{\theta})$  (the second part can be ignored)
    - M-step objective is upper bounded by true objective, and is equal to true objective at current parameter estimate. M-step optimization can be done efficiently in most cases
- A example:



- $L = l$  at  $\theta_t$  because of the E-step
- EM Convergence
  - Assuming  $\boldsymbol{\theta}^{(t)}$  and  $\boldsymbol{\theta}^{(t+1)}$  are the parameters from two successive iterations of EM, we have

$$\begin{aligned}
 & L(\boldsymbol{\theta}^{(t+1)}) \\
 &= l(\boldsymbol{\theta}^{(t)}) = \sum_{i=1}^m \sum_{z^{(i)} \in Z} q(z^{(i)}; \boldsymbol{\theta}^{(t)}) \log \frac{p(\mathbf{x}^{(i)}, z^{(i)}; \boldsymbol{\theta}^{(t)})}{q(z^{(i)}; \boldsymbol{\theta}^{(t)})} \\
 &\leq \sum_{i=1}^m \sum_{z^{(i)} \in Z} q(z^{(i)}; \boldsymbol{\theta}^{(t)}) \log \frac{p(\mathbf{x}^{(i)}, z^{(i)}; \boldsymbol{\theta}^{(t+1)})}{q(z^{(i)}; \boldsymbol{\theta}^{(t)})} \\
 &\leq l(\boldsymbol{\theta}^{(t+1)}) \\
 &= \sum_{i=1}^m \sum_{z^{(i)} \in Z} q(z^{(i)}; \boldsymbol{\theta}^{(t+1)}) \log \frac{p(\mathbf{x}^{(i)}, z^{(i)}; \boldsymbol{\theta}^{(t+1)})}{q(z^{(i)}; \boldsymbol{\theta}^{(t+1)})} \\
 &= L(\boldsymbol{\theta}^{(t+1)})
 \end{aligned}$$

- Hence, EM causes the likelihood to increase monotonically, and if bounded above, converges to some  $L^*$
- Remark: If we define:

$$J(q, \boldsymbol{\theta}) = \sum_{i=1}^m \sum_{z^{(i)} \in Z} q(z^{(i)}) \log \frac{p(\mathbf{x}^{(i)}, z^{(i)}; \boldsymbol{\theta})}{q(z^{(i)})},$$

EM can also be viewed coordinate ascent on  $J$ , in which the E-step maximizes  $J$  w.r.t.  $q$ , and the M-step maximizes  $J$  with respect to  $\boldsymbol{\theta}$

- Proof:
  - For the objective function:

$$J(q, \theta) = \sum_{i=1}^m \sum_{z^{(i)} \in Z} q(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{q(z^{(i)})}$$

where  $\sum_{z^{(i)} \in Z} q(z^{(i)}) = 1$

- We want to maximize  $J$ , which is a **concave function**, we reverse the Lagrange dual form:

- Lagrangian is :

$$\begin{aligned} L(q, \theta) &= J(q, \theta) + \lambda \left( 1 - \sum_{z^{(i)} \in Z} q(z^{(i)}) \right) \\ &= \sum_{i=1}^m \sum_{z^{(i)} \in Z} q(z^{(i)}) \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{q(z^{(i)})} + \lambda \left( 1 - \sum_{z^{(i)} \in Z} q(z^{(i)}) \right) \end{aligned}$$

- Take the partial derivative of  $L$  w.r.t  $q(z^{(i)})$  and setting it to zero, we have:

$$\begin{aligned} \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{q(z^{(i)})} + q(z^{(i)}) \frac{q(z^{(i)})}{p(x^{(i)}, z^{(i)}; \theta)} \left( -\frac{p(x^{(i)}, z^{(i)}; \theta)}{q(z^{(i)})^2} \right) - \lambda &= 0 \\ \log \frac{p(x^{(i)}, z^{(i)}; \theta)}{q(z^{(i)})} &= 1 + \lambda \\ g(z^{(i)}) &= \frac{p(x^{(i)}, z^{(i)}; \theta)}{q(z^{(i)})} = e^{1+\lambda} (\text{a constant}) \end{aligned}$$

- Therefore, when we maximize  $J$  w.r.t  $q$ , we fulfill the condition that  $g$  is a constant, which make sure once we use  $q(z^{(i)}; \theta^{(t+1)})$ , we maximize the function  $J$  which is an equivalent to maximizing  $l(q, \theta)$ . At the same time  $l$  is equal to  $L$ , which means we are maximizing  $L$

$$\begin{aligned} L(\theta^{(t+1)}, q(z^i, \theta^{(t+1)})) &> L(\theta^{(t+1)}, q(z^i, \theta^{(t)})) \\ &> l(\theta^{(t+1)}, q(z^i, \theta^{(t)})) > l(\theta^{(t)}, q(z^i, \theta^{(t)})) = L(\theta^{(t)}, q(z^i, \theta^{(t)})) \end{aligned}$$

- Continue calculating the  $q(z^{(i)})$

$$\begin{aligned} \sum_{z^{(i)} \in Z} q(z^{(i)}) &= \sum_{z^{(i)} \in Z} \frac{p(x^{(i)}, z^{(i)}; \theta)}{e^{1+\lambda}} = \frac{p(x^{(i)}; \theta)}{e^{1+\lambda}} = 1 \\ q(z^{(i)}) &= \frac{p(x^{(i)}, z^{(i)}; \theta)}{e^{1+\lambda}} = p(z^{(i)} | x^{(i)}; \theta) \end{aligned}$$

- Which is exactly the M-step

## Gaussian Mixture Models

- Mixture Models
  - A mixture model is a probabilistic clustering model that is the unsupervised analogue(模拟, 近似) of the Bayes optimal classifier where the **unknown assignment of data cases** to clusters take the place of the known class labels
  - We let  $x^{(i)}$  be a data case and  $z^{(i)} \in \{1, \dots, K\}$  be the index of the cluster  $x^{(i)}$  belongs to.  $z^{(i)}$  is often called the mixture indicator variable or the latent class (i.e. unobserved)
  - Each cluster  $k$  specifies its **own distribution over the feature vectors**  $p(x^{(i)} | z^{(i)} = k)$
  - We also have a discrete distribution  $p(z^{(i)} = k) = \theta_k$ , which describes the prior probability that  $x^{(i)}$  belongs to cluster  $k$
- Data Distribution
  - The joint distribution of the feature( $x^{(i)}$ ) and the mixture indicator variable( $z^{(i)}$ ) is :

$$p(\mathbf{x}^{(i)}, z^{(i)} = k) = p(\mathbf{x}^{(i)} | z^{(i)} = k)p(z^{(i)} = k)$$

- In clustering, we don't know that the right value of the mixture indicator variable is a priori, but we can marginalize it away to obtain a probability distribution on the feature vector only

$$p(\mathbf{x}^{(i)}) = \sum_{k=1}^K p(\mathbf{x}^{(i)} | z^{(i)} = k)p(z^{(i)} = k)$$

- Mixture Component Distributions: Common model choices for  $p(\mathbf{x}|z = k)$  are:

- Categorical:  $\prod_{j=1}^n \prod_{x \in X_j} \theta_{x,j|k}^{\mathbb{I}[x_j=x]}$ , where  $X_j$  contains all possibilities of  $j^{th}$  feature of  $\mathbf{x}$
- Independent Gaussian:  $\prod_{j=1}^n N(x_j; \mu_{j|k}, \sigma_{j|k}^2)$
- Multivariate Gaussian:  $N(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$

- Learning:

- Given data set  $D = \{\mathbf{x}^{(i)}\}_{i=1,\dots,m}$ , we can learn the mixture model parameters by maximizing the log probability of the data give the parameters:

$$l = \sum_{i=1}^m \log\left(\sum_{k=1}^K p(\mathbf{x}^{(i)} | z^{(i)} = k)p(z^{(i)} = k)\right)$$

- While we can do this directly using gradient-based optimization, it's often faster to use a special algorithm called *Expectation Maximization*

- Expectation Maximization for Gaussian Mixture Models

- E-Step: In the first step of the algorithm, we compute the posterior probability that each data case belongs to each cluster using the Bayes rule:

$$w_k^{(i)} = p(z^{(i)} = k | \mathbf{x}^{(i)}) = \frac{\theta_k N(\mathbf{x}^{(i)}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{k'=1}^K \theta_{k'} N(\mathbf{x}^{(i)}; \boldsymbol{\mu}_{k'}, \boldsymbol{\Sigma}_{k'})}$$

- M-Step: In the second step, we update the parameters using the weighted averages:

$$\theta_k = \frac{\sum_{i=1}^m w_k^{(i)}}{m}, \quad \boldsymbol{\mu}_k = \frac{\sum_{i=1}^m w_k^{(i)} \mathbf{x}^{(i)}}{\sum_{i=1}^m w_k^{(i)}}$$

$$\boldsymbol{\Sigma}_k = \frac{\sum_{i=1}^m w_k^{(i)} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_k)(\mathbf{x}^{(i)} - \boldsymbol{\mu}_k)^T}{\sum_{i=1}^m w_k^{(i)}}$$

- A Special Case

- Suppose we fix  $\theta_k = \frac{1}{K}$  and  $\sum_k = \mathbf{I}$ . In this case we have:

$$N(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \frac{1}{\sqrt{|(2\pi)^n I|}} \exp\left(-\frac{1}{2} \|\mathbf{x}^{(i)} - \boldsymbol{\mu}_k\|_2^2\right)$$

- and we obtain the following special case of the EM algorithm for multivariate Gaussians:

$$w_k^{(i)} = \frac{\exp(-\frac{1}{2} \|\mathbf{x}^{(i)} - \boldsymbol{\mu}_k\|_2^2)}{\sum_{k'=1}^K \exp(-\frac{1}{2} \|\mathbf{x}^{(i)} - \boldsymbol{\mu}_{k'}\|_2^2)},$$

$$\boldsymbol{\mu}_k = \frac{\sum_{i=1}^m w_k^{(i)} \mathbf{x}^{(i)}}{\sum_{i=1}^m w_k^{(i)}}$$

This is often referred to as soft  $K$ -means.

---

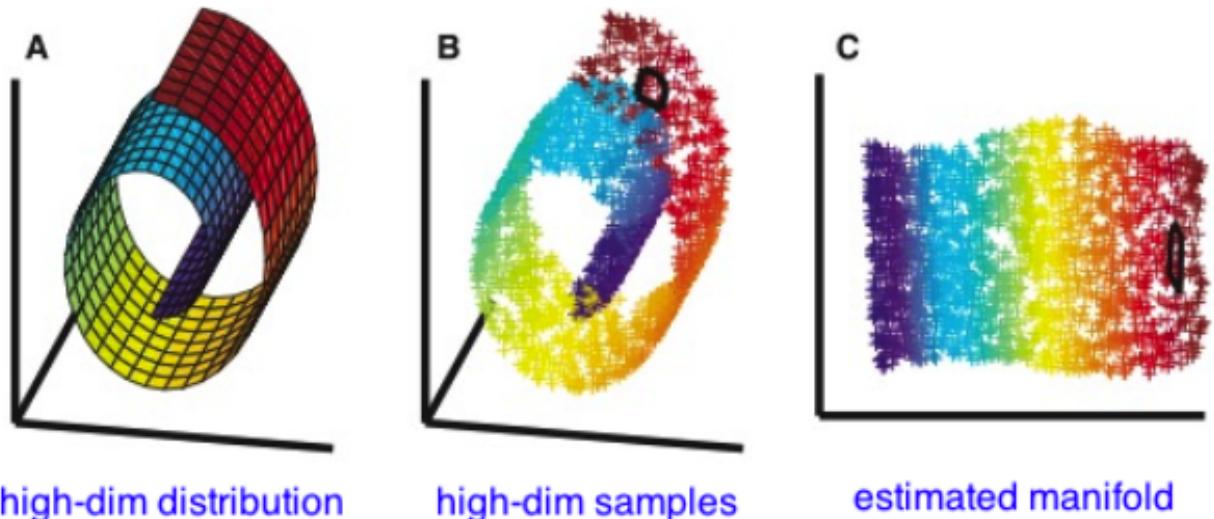
- Trade-Offs

- We can see that the original  $K$ -means algorithm performs **hard assignments during clustering**, and implicitly assumes all clusters will have an equal number of points assigned as well as a unit covariance matrix
- EM for Gaussian mixture models relaxes all of these assumptions. The objective still has multiple local optima, but **EM also produces a guaranteed non-decreasing sequence of objective function values**
- EM can also be used with any component densities/distributions to customize the model to a given data set
- As with  $K$ -means, initialization is important, but the same heuristics can be applied

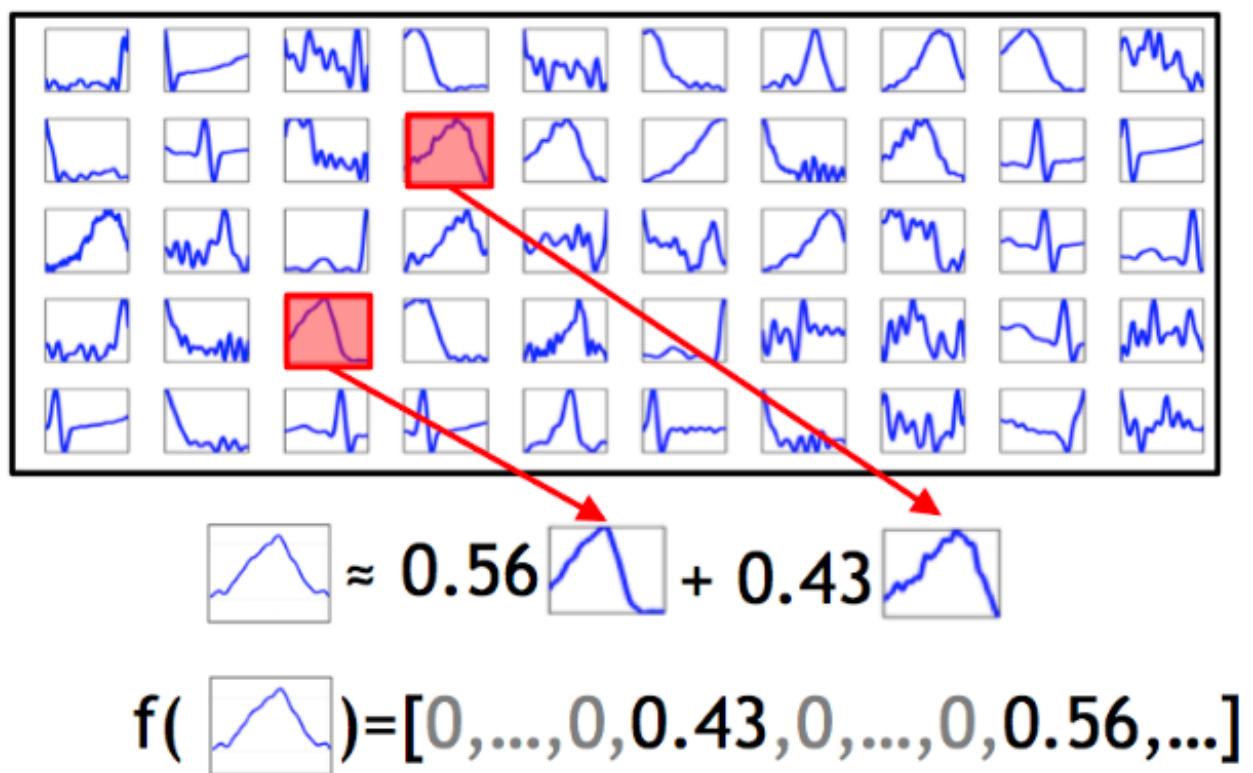
# Lecture 09 Linear Dimensionality Reduction

## Dimensionality Reduction

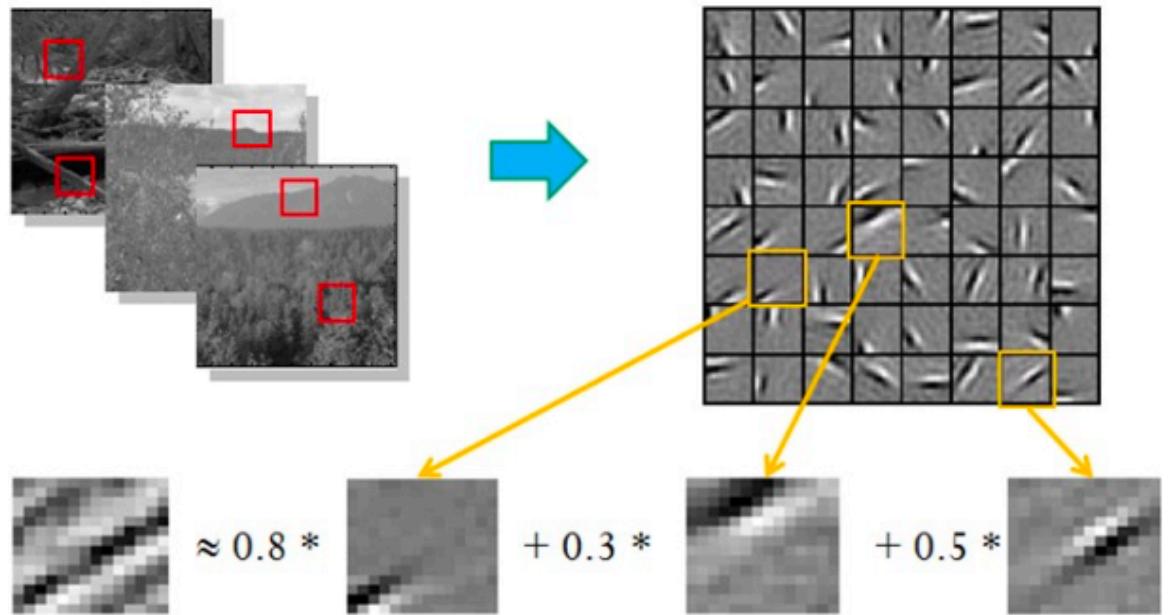
- Definition of dimensionality reduction task:
  - Given a collection of feature vectors  $\{\mathbf{x}^{(i)} \in \mathbb{R}^n\}$ , map the feature vectors into a lower dimensional space  $\{\mathbf{z}^{(i)} \in \mathbb{R}^l\}$ , where  $l < n$ , while preserving certain “properties” of the data



- Example
  - Representing Time Series



- Representing Natural Image Patches



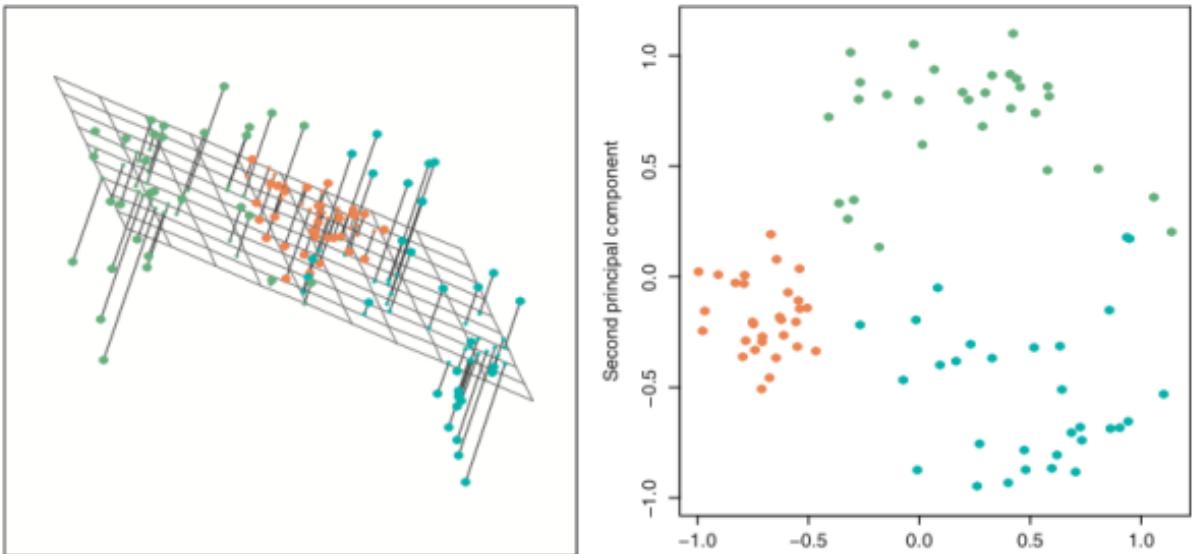
$$[\mathbf{a}_1, \dots, \mathbf{a}_{64}] = [0, 0, \dots, 0, \mathbf{0.8}, 0, \dots, 0, \mathbf{0.3}, 0, \dots, 0, \mathbf{0.5}, 0]$$

(feature representation)

- Application
  - Can be used as a **pre-processing step** to enable more accurate classification and regression on manifold data
  - Very low dimensional embeddings (i.e.,  $l = 2, 3$ ) can be used to **visualize complex manifold-structured data** as in the previous example
  - Can be used to “**de-noise**” data by **projecting to lower-dimensional space** and then **projecting back to the feature space**
- Dimensionality Reduction vs. Feature Selection
  - The goal of feature selection is to **remove features that are not informative** with respect to the class label. This obviously reduces the dimensionality of the feature space
  - Dimensionality reduction can be used to find a **meaningful lower dimensional feature space** for manifold-structured data even when there is information in each of the feature dimensions so that none can be discarded (Linear independent)
  - Another important property of dimensionality reduction is that it is **unsupervised**. It will attempt to **preserve(保留) structure in the data that could be useful for a range of supervised problems**
  - Unlike feature selection, which is a supervised task, dimensionality reduction can sometimes **discard structure useful to a particular supervised task** thereby increasing error
- Dimensionality Reduction vs. Data Compression
  - While dimensionality reduction can be seen as a simplistic form of data compression, it is not equivalent to it, as the goal of data compression is to reduce the **entropy** of the representation **not only** the dimensionality
  - For example, in lossless data compression, *arithmetic coding* encodes the entire data into a single number, an arbitrary-precision fraction  $q$  where  $0.0 \leq q < 1.0$ . In some science fiction books, this is paraphrased as a pinpoint representing all information of the whole universe

# Linear Dimensionality Reduction

- Linear Dimensionality Reduction:
  - The simplest dimensionality reduction methods assumes that the observed high dimensional data vectors  $\mathbf{x}^{(i)} \in \mathbb{R}^n$  lie on an  $l$ -dimensional linear subspace with  $\mathbb{R}^n$  (Project the data vectors to a lower dimensional space)



- Assumption: Firstly, assume all the data points are already on a lower dimensional space
    - Mathematically, the linear subspace assumption can be written as follows:
- $$\mathbf{x}^{(i)} = \sum_{k=1}^l z_k^{(i)} \mathbf{b}_k \quad (1)$$
- Where  $\mathbf{b}_k = [b_{k1}, b_{k2}, \dots, b_{kn}]^T$  for  $b_{ki}$  where  $i = 1, \dots, l$  are a set of basis vectors describing an  $l$ -dimensional linear sub-space of  $\mathbb{R}^n$  and  $z_k^{(i)}$ 's are real-valued weights on those basis vectors (reduced data vectors)
- Connection to linear regression
    - This expression is exactly linear regression where  $x_j^{(i)}$  is the target,  $z_k^{(i)}$ 's are the weights, and  $b_{kj}$  for each  $k$  are the features
- $$x_j^{(i)} = \sum_{k=1}^l z_k^{(i)} b_{kj} \quad (2)$$
- This observation is also true if we swap the roles of the weights and the features
  - However, unlike linear regression, we only know what corresponds to the targets. We must learn both the features and the weights
- Matrix Form
    - If we let  $\mathbf{X}$  be the data matrix  $X_{ij} = x_j^{(i)}$ ,  $\mathbf{Z}$  be a matrix where  $Z_{ik} = z_k^{(i)}$ , and  $\mathbf{B}$  be a matrix where  $B_{kj} = b_{kj}$ , we can express  $\mathbf{X}$  as follows:

$$\mathbf{X} = \mathbf{Z} \times \mathbf{B} \quad (3)$$

- $\mathbf{Z} \in \mathbb{R}^{m \times l}$  is often referred to as the **factor loading matrix** while  $\mathbf{B} \in \mathbb{R}^{l \times n}$  are referred to as the **latent factors**
- Most real world data will be subject to noise. If we assume that  $\epsilon \in \mathbb{R}^{m \times n}$  is a matrix of noise values from some probability distribution, we have:

$$\mathbf{X} = \mathbf{Z} \times \mathbf{B} + \epsilon \quad (4)$$

- Learning

- The learning problem for linear dimensionality reduction is to estimate values for both  $Z$  and  $B$  given only the noisy observations  $X$
- One possible learning criteria is to minimize the sum of squared errors when reconstructing  $\mathbf{X}$  from  $\mathbf{Z}$  and  $\mathbf{B}$ . This leads to:

$$\operatorname{argmin}_{\mathbf{Z}, \mathbf{B}} \|\mathbf{X} - \mathbf{Z}\mathbf{B}\|_F^2 \quad (5)$$

- where  $\|\mathbf{A}\|_F$  is the Frobenius norm of matrix  $\mathbf{A}$  ((the square root of the sum of the squares of all matrix entries => Least Square)

- Some properties of **trace** in a matrix and the Frobenius norm:

We first recall the definition of the *trace* of a square matrix  $A \in \mathbb{R}^{n \times n}$

$$\operatorname{tr} A = \sum_{i=1}^n A_{ii} = A_{11} + A_{22} + \dots + A_{nn}. \quad (135)$$

The trace is a linear mapping. That is,

$$\operatorname{tr}(A + B) = \operatorname{tr} A + \operatorname{tr} B, \quad (136)$$

and

$$\text{tr}cA = c\text{tr}A. \quad (137)$$

We note a trivial but useful property of the trace

$$\text{tr}A = \text{tr}A^T = \sum_{j=1}^n A_{jj}. \quad (138)$$

This result implies that for  $A \in \mathbb{R}^{m \times n}$  and  $B \in \mathbb{R}^{m \times n}$

$$\text{tr}AB^T = \text{tr}(AB^T)^T = \text{tr}BA^T \quad (139)$$

and

$$\text{tr}A^TB = \text{tr}(A^TB)^T = \text{tr}B^TA. \quad (140)$$

We prove another important property of the trace: for  $A \in \mathbb{R}^{m \times n}$  and  $B \in \mathbb{R}^{n \times m}$ ,  $\text{tr}AB = \text{tr}BA$ .

*Proof.*

$$\text{tr}AB = \sum_{i=1}^m (AB)_{ii} = \sum_{i=1}^m \sum_{j=1}^n a_{ij}b_{ji} = \sum_{j=1}^n \sum_{i=1}^m b_{ji}a_{ij} = \sum_{j=1}^n (BA)_{jj} = \text{tr}BA. \quad (141)$$

□

More generally, the trace is invariant under cyclic permutations, which can be easily proved using Eq. (141). For example, for  $A \in \mathbb{R}^{m \times n}$ ,  $B \in \mathbb{R}^{n \times l}$ , and  $C \in \mathbb{R}^{l \times m}$ , we have

$$\text{tr}ABC = \text{tr}A(BC) = \text{tr}BCA = \text{tr}B(CA) = \text{tr}CAB. \quad (142)$$

We now prove that for  $A \in \mathbb{R}^{m \times n}$

---


$$\|A\|_F^2 = \text{tr}AA^T = \text{tr}A^TA. \quad (143)$$

*Proof.*

$$\|A\|_F^2 = \sum_{i=1}^m \sum_{j=1}^n A_{ij}^2 = \sum_{i=1}^m \left( \sum_{j=1}^n A_{ij}A_{ij} \right) = \sum_{i=1}^m \left( \sum_{j=1}^n A_{ij}A_{ji}^T \right) \quad (144)$$

$$= \sum_{i=1}^m (AA^T)_{ii} = \text{tr}AA^T. \quad (145)$$

In Eq. (144), we note the fact that the entry in the  $i$ th row and  $j$ th column of  $A$  is equal to the entry in the  $j$ th row and  $i$ th column of  $A^T$ . □

Using Eq. (143), we are able to decompose the loss function in linear dimensionality reduction as

$$\begin{aligned}\ell(Z, B) &= \|X - ZB\|_F^2 = \text{tr}(X - ZB)(X - ZB)^T \\ &= \text{tr}(XX^T - X(ZB)^T - ZBX^T + ZB(ZB)^T) \\ &= \text{tr}(XX^T - 2ZBX^T + ZBB^TZ^T) \quad (146) \\ &= \text{tr}XX^T - 2\text{tr}ZBX^T + \text{tr}ZBB^TZ^T. \quad (147)\end{aligned}$$

In Eq. (146), we make use of the trace identity  $\text{tr}AB^T = \text{tr}BA^T$ . To derive the closed-form update rules of alternating least squares in solving the linear dimensionality reduction problem, we first take the derivative of  $\ell$  w.r.t.  $Z$

$$\begin{aligned}\frac{\partial \ell(Z, B)}{\partial Z} &= -2\frac{\partial}{\partial Z}\text{tr}Z(BX^T) + \frac{\partial}{\partial Z}\text{tr}Z(BB^TZ^T) \\ &= -2XB^T + ZBB^T + Z(BB^T)^T \\ &= -2XB^T + 2ZBB^T, \quad (148)\end{aligned}$$

where we make use of two derivatives of traces<sup>22</sup>

$$\frac{\partial}{\partial X}\text{tr}XA = A^T \quad (149)$$

and

$$\frac{\partial}{\partial X}\text{tr}XAX^T = XA + XA^T. \quad (150)$$

Setting the partial derivative to zero, we have

$$ZBB^T = XB^T \quad (151)$$

$$Z = XB^T(BB^T)^{-1} \quad (152)$$

$$Z^T = (XB^T(BB^T)^{-1})^T \quad (153)$$

$$= (BB^T)^{-1}BX^T. \quad (154)$$

In Eq. (154), we exploit the fact that the inverse of the transpose of an invertible matrix is equal to the transpose of its inverse, *i.e.*,  $(A^T)^{-1} = (A^{-1})^T$ . Similarly, taking the derivative of  $\ell$  w.r.t.  $B$ , we have

$$\begin{aligned}\frac{\partial \ell(Z, B)}{\partial B} &= -2\frac{\partial}{\partial B}\text{tr}ZBX^T + \frac{\partial}{\partial B}\text{tr}ZBB^TZ^T \\ &= -2\frac{\partial}{\partial B}\text{tr}B(X^TZ) + \frac{\partial}{\partial B}\text{tr}(Z^TZ)BB^T \\ &= -2Z^TX + Z^TZB + (Z^TZ)^TB \\ &= -2Z^TX + 2Z^TZB, \quad (155)\end{aligned}$$

where we use the derivatives of trace<sup>23</sup>

$$\frac{\partial}{\partial X} \text{tr} AXX^T = AX + A^T X. \quad (156)$$

Setting the partial derivative to zero, we have

$$\begin{aligned} Z^T Z B &= Z^T X \\ B &= (Z^T Z)^{-1} Z^T X. \end{aligned} \quad (157)$$

- Learning Algorithm

- Not surprisingly, we can obtain a solution to this learning problem by leveraging the OLS solution to linear regression. The algorithm is often referred to as Alternating Least Squares (ALS)
- Starting from a random initialization, ALS iterates between **assuming  $Z$  are known features and optimizing  $B$  as the unknown weights**, and **assuming that  $B$  are the known features and optimizing  $Z$  as the unknown weights** (is it coordinate descent?)

$$\begin{aligned} \mathbf{B} &\leftarrow (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{X}, \\ \mathbf{Z}^T &\leftarrow (\mathbf{B} \mathbf{B}^T)^{-1} \mathbf{B} \mathbf{X}^T \end{aligned}$$

- Lack of Uniqueness of Optimal Parameters

- Suppose we run the ALS algorithm to convergence and obtain estimates for  $\mathbf{Z}^*$  and  $\mathbf{B}^*$  s.t.:

$$l^* = \|\mathbf{X} - \mathbf{Z}^* \mathbf{B}^*\|_F^2 \quad (6)$$

- Note that if we let  $\mathbf{R} \in \mathbb{R}^{l \times l}$  be an arbitrary  $l \times l$  invertible matrix, then we obtain exactly the same value  $l^*$  of the objective function for the alternate parameters  $Z' = Z^* R$  and  $B' = R^{-1} B^*$

$$\begin{aligned} l^* &= \|\mathbf{X} - \mathbf{Z}^* (\mathbf{I}) \mathbf{B}^*\|_F^2 \\ &= \|\mathbf{X} - \mathbf{Z}^* (\mathbf{R} \mathbf{R}^{-1}) \mathbf{B}^*\|_F^2 \\ &= \|\mathbf{X} - \mathbf{Z}' \mathbf{B}'\|_F^2 \end{aligned} \quad (7)$$

## Singular Value Decomposition

- SVD: ([resource](#))

- Let  $\mathbf{X}$  be an  $m \times n$  matrix, with  $m \geq n$ . It can be factorized

$$\mathbf{X} = \mathbf{U} \begin{pmatrix} \Sigma \\ 0 \end{pmatrix} \mathbf{V}^T \quad (8)$$

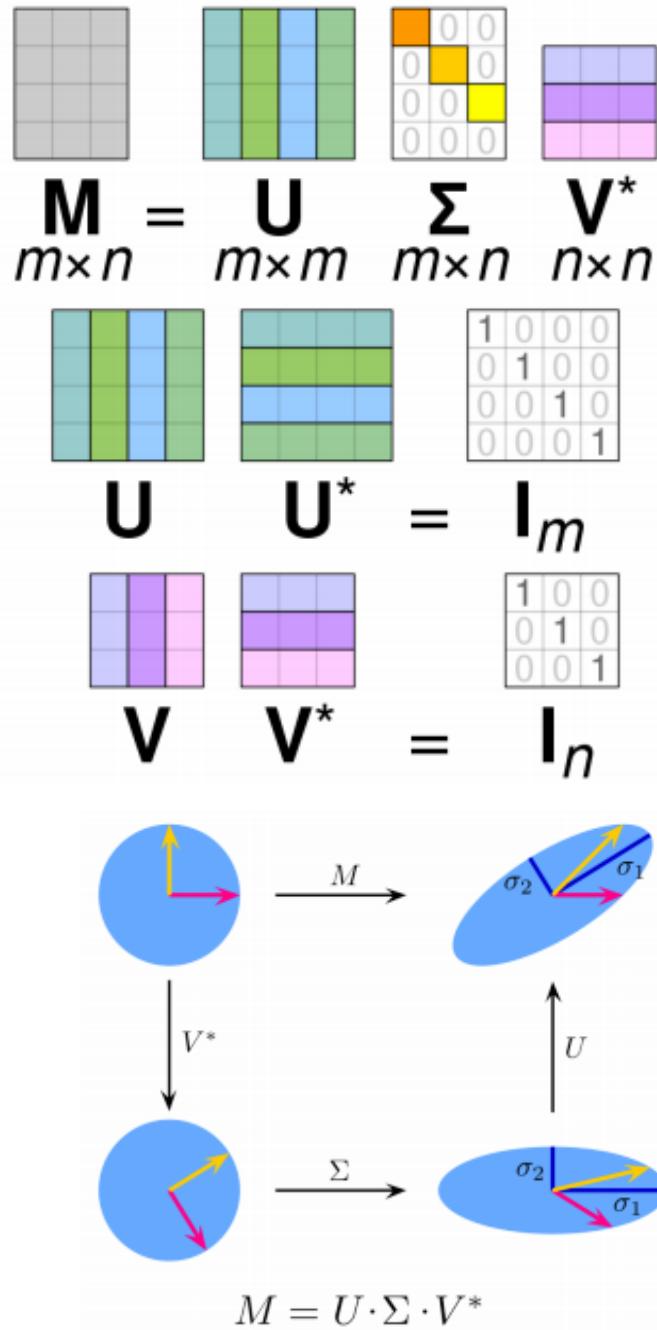
- where  $\mathbf{U} \in \mathbb{R}^{m \times m}$  and  $\mathbf{V} \in \mathbb{R}^{n \times n}$  are orthogonal, i.e.,

$$\begin{aligned} \mathbf{U}^T \mathbf{U} &= \mathbf{U} \mathbf{U}^T = \mathbf{I}_m \\ \mathbf{V}^T \mathbf{V} &= \mathbf{V} \mathbf{V}^T = \mathbf{I}_n \end{aligned} \quad (9)$$

- and  $\Sigma \in \mathbb{R}^{n \times n}$  is diagonal (This matrix rescale the vectors in  $V$ )

$$\Sigma = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_n), \sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n \geq 0 \quad (10)$$

- The quantities  $\sigma_i$ 's are called *singular(eigen) values* of  $\mathbf{X}$  and the columns of  $\mathbf{U}$  and  $\mathbf{V}$  are called the left and right *singular(eigen) vectors* of  $\mathbf{X}$ , respectively



- Upper left: The unit disc with the two canonical unit vectors
- Upper right: transformed with  $\mathbf{M}$
- Lower left: The action of  $\mathbf{V}^T$ . This is just a rotation
- Lower right: The action of  $\Sigma \mathbf{V}^T$ .  $\Sigma$  scales in vertically and horizontally

- Properties of each matrix

- Matrix  $\mathbf{U}$

- Columns of  $\mathbf{U}$  consist of eigenvectors of  $\mathbf{XX}^T$

$$\begin{aligned}\mathbf{XX}^T &= (\mathbf{U}\Sigma\mathbf{V}^T)(\mathbf{U}\Sigma\mathbf{V}^T)^T \\ &= (\mathbf{U}\Sigma\mathbf{V}^T)(\mathbf{V}\Sigma^T\mathbf{U}^T) \\ &= \mathbf{U}\Sigma\Sigma^T\mathbf{U}^T = \mathbf{U}(\Sigma)^2\mathbf{U}^T\end{aligned}\tag{11}$$

- Orthogonal and has unit norm, *i.e.*,  $\mathbf{UU}^T = \mathbf{U}^T\mathbf{U} = \mathbf{I}$

- Matrix  $\Sigma$

- Diagonal
- Square roots of the eigenvalues of both  $\mathbf{X}^T\mathbf{X}$  and  $\mathbf{XX}^T$
- Shows importance of each eigenvector

- Matrix  $\mathbf{V}$

- Columns of  $\mathbf{V}$  consist of eigenvectors of  $\mathbf{X}^T\mathbf{X}$

$$\begin{aligned}\mathbf{X}^T\mathbf{X} &= (\mathbf{U}\Sigma\mathbf{V}^T)^T(\mathbf{U}\Sigma\mathbf{V}^T) \\ &= (\mathbf{V}\Sigma^T\mathbf{U}^T)(\mathbf{U}\Sigma\mathbf{V}^T) \\ &= \mathbf{V}\Sigma^T\Sigma\mathbf{V}^T = \mathbf{V}(\Sigma)^2\mathbf{V}^T\end{aligned}\tag{12}$$

- Orthogonal and has unit form. *i.e.*,  $\mathbf{VV}^T = \mathbf{V}^T\mathbf{V} = \mathbf{I}$

- Reduced-Form Singular Value Decomposition (SVD)

- If only  $l < \min\{m, n\}$  singular values are non-zeros, the SVD of  $\mathbf{X} \in \mathbb{R}^{m \times n}$  can be represented in *reduced form* as follows:

$$\mathbf{X} = \mathbf{U}\Sigma_l\mathbf{V}^T = \sum_{j=1}^l \sigma_j \mathbf{u}_j \mathbf{v}_j^T\tag{13}$$

- where  $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_l]$  and  $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_l]$ , satisfying

$$\mathbf{U}^T\mathbf{U} = \mathbf{I}_l, \quad \mathbf{V}^T\mathbf{V} = \mathbf{I}_l\tag{14}$$

- and  $\Sigma_l = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_l)$ ,  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_l \geq 0$
- $\mathbf{u}_j \mathbf{v}_j^T \in \mathbb{R}^{m \times n}$  is the product of a column vector  $\mathbf{u}_j$  and the transpose of a column vector  $\mathbf{v}_j$ . It has rank 1. That is,  $\mathbf{X}$  is a weighted summation of  $l$  rank-1 matrices in  $\mathbb{R}^{m \times n}$

- $\Sigma_l$  stores a combination of  $l$  rows in  $\mathbf{V}$  (*it scales* each vectors in  $\mathbf{V}$ )

- When  $m > n$ , it works like:

- $\mathbf{U}$  is a matrice in  $\mathbb{R}^{m \times m}$ , and the  $l$  column vectors  $\mathbf{u}^{(j)}$  corresponds to the  $\sigma_i$  are moved to the corresponding places
- $\mathbf{V}$  is a matrice in  $\mathbb{R}^{n \times n}$ , and the  $l$  row vectors  $\mathbf{v}^{(i)}$  corresponds to the  $\sigma_i$  are moved to the corresponding places
- $\Sigma_l$  and  $\mathbf{V}'$  are matrices in  $\mathbb{R}^{l \times l}$
- $(\mathbf{V}' \ 0)$  and  $(\Sigma_l \ 0)$  are matrices in  $\mathbb{R}^{m \times n}$

$$\mathbf{U} \begin{pmatrix} (\mathbf{V}')^T & 0 \\ 0 & 0 \end{pmatrix} = \mathbf{U} \begin{pmatrix} \Sigma_l & 0 \\ 0 & 0 \end{pmatrix} \mathbf{V}^T\tag{15}$$

- When  $l = n$

$$\mathbf{U} \begin{pmatrix} (\mathbf{V}')^T \\ 0 \end{pmatrix} = \mathbf{U} \begin{pmatrix} \boldsymbol{\Sigma} \\ 0 \end{pmatrix} \mathbf{V}^T \quad (16)$$

## Eckart-Young-Mirsky Theorem

- Given an  $m \times n$  matrix  $\mathbf{X}$  of rank  $r \leq \min\{m, n\}$  ([resource1](#) [resource2](#)) and its singular value decomposition  $\mathbf{X} = \mathbf{U}\boldsymbol{\Sigma}_r\mathbf{V}^T$ , with singular values  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r$ , then among all  $m \times n$  matrices of lower rank  $l < r$ , the best approximation is  $\mathbf{Y}^* = \mathbf{U}\boldsymbol{\Sigma}_l\mathbf{V}^T$ , where  $\boldsymbol{\Sigma}_l$  is the **diagonal matrix with singular values  $\sigma_1, \sigma_2, \dots, \sigma_l$  in the sense that**

$$\|\mathbf{X} - \mathbf{Y}^*\|_F^2 = \min\{\|\mathbf{X} - \mathbf{Y}\|_F^2; \mathbf{Y} \in \mathbb{R}^{m \times n}, \text{rank } \mathbf{Y} \leq l\} \quad (17)$$

- This is not a convex optimization problem since the condition function  $\text{rank}$  is **neither convex nor concave**:

$$X_1 = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}, X_2 = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} \quad (18)$$

$$\text{rank}(X_1) = \text{rank}(X_2) = 1$$

$$Y = \theta X_1 + (1 - \theta) X_2 = \begin{bmatrix} \theta & 0 \\ 0 & (1 - \theta) \end{bmatrix}, \text{where } \theta \in [0, 1]$$

$$\text{rank}(Y) = 2, \text{when } \theta \in (0, 1)$$

$$2 = \text{rank}(\theta X_1 + (1 - \theta) X_2) > \theta \text{rank}(X_1) + (1 - \theta) \text{rank}(X_2) = 1$$

$$X_1 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, X_2 = \begin{bmatrix} 0 & 0 \\ 0 & -1 \end{bmatrix}$$

$$\text{rank}(X_1) = 2, \text{rank}(X_2) = 1$$

$$Y = \theta X_1 + (1 - \theta) X_2 = \begin{bmatrix} \theta & 0 \\ 0 & 2\theta - 1 \end{bmatrix}, \text{where } \theta \in [0, 1]$$

$$\text{rank}(\mathbf{Y}) = 1, \text{when } \theta = \frac{1}{2}$$

$$1 = \text{rank}(\theta X_1 + (1 - \theta) X_2) < \theta \text{rank}(X_1) + (1 - \theta) \text{rank}(X_2) = \frac{3}{2}$$

- Proof: Since the Frobenius norm is preserved under unitary transformations, we have

$$\|\mathbf{X} - \mathbf{Y}\|_F^2 = \|\mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T - \mathbf{Y}\|_F^2 = \|\boldsymbol{\Sigma} - \mathbf{U}^T\mathbf{Y}\mathbf{V}\|_F^2 \quad (19)$$

(for orthogonal matrix  $\mathbf{U}$ )

$$\|\mathbf{U}\mathbf{X}\|_F^2 = \text{tr}((\mathbf{U}\mathbf{X})^T(\mathbf{U}\mathbf{X})) = \text{tr}(\mathbf{X}^T\mathbf{U}^T\mathbf{U}\mathbf{X})$$

$$= \text{tr}(\mathbf{X}^T\mathbf{X}) = \|\mathbf{X}\|_F^2$$

$$\|\mathbf{X}\mathbf{U}\|_F^2 = \text{tr}((\mathbf{X}\mathbf{U})^T(\mathbf{X}\mathbf{U})) = \text{tr}(\mathbf{U}^T\mathbf{X}^T\mathbf{X}\mathbf{U})$$

$$= \text{tr}(\mathbf{X}^T\mathbf{X}\mathbf{U}\mathbf{U}^T) = \text{tr}(\mathbf{X}^T\mathbf{X}) = \|\mathbf{X}\|_F^2$$

Denoting  $\mathbf{Z} = \mathbf{U}^T\mathbf{Y}\mathbf{V}$ , an  $m \times n$  matrix of rank  $l$ , a direct calculation gives

$$\|\boldsymbol{\Sigma} - \mathbf{Z}\|_F^2 = \sum_{ij} |\Sigma_{ij} - Z_{ij}|^2 = \sum_{i=1}^r |\sigma_i - Z_{ij}|^2 + \sum_{i>r} |Z_{ii}|^2 (= 0) + \sum_{i \neq j} |Z_{ij}|^2$$

- If  $\mathbf{Z}$  is not diagonal, set  $Z_{ij} = 0 (i \neq j)$  will have a better result
- if we don't use exactly the largest (mean value)  $\sigma_1, \sigma_2, \dots, \sigma_l$  to match the larger  $l \sigma$  in

$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r$ , there will also be some better result

- Therefore,  $\mathbf{Z}^* = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_l) = \boldsymbol{\Sigma}_l$ ;  $\mathbf{Y}^* = \mathbf{U}\boldsymbol{\Sigma}_l\mathbf{V}^T$

- Trade-Offs

- Minimum Frobenius norm linear dimensionality reduction is nice because it will find the globally optimal linear sub-space
- The basic ALS(Alternating Least Square) algorithm scales as  $O(l^3 + ml^2 + nl^2 + mnl)$  per pair of iterations. The full SVD algorithm scales as  $O(\min(mn^2, nm^2))$ . Fast approximations exist for both problems
- A significant limitation of linear dimensionality reduction is that it can fail to achieve a useful compact(dense) representation of the data if the underlying manifold is not actually linear (non-unique)
- As with clustering, the rank  $l$  of the latent sub-space is a free parameter

## Principal Component Analysis

- Eigenvectors

- Assume  $\mathbf{A} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{v} \in \mathbb{C}^{n \times 1}$ , and  $\lambda \in \mathbb{C}$
- If  $\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$  the  $\mathbf{v}$  is a right eigenvector of  $\mathbf{A}$  with eigenvalue  $\lambda$
- If  $\mathbf{A}^T\mathbf{v} = \lambda\mathbf{v}$  ( $\mathbf{v}^T\mathbf{A} = \lambda\mathbf{v}^T$ ) the  $\mathbf{v}$  is a right eigenvector of  $\mathbf{A}$  with eigenvalue  $\lambda$
- If  $\mathbf{A}$  is symmetric so that  $\mathbf{A} = \mathbf{A}^T$ , then the left and right eigenvectors of  $\mathbf{A}$  are the same with the same eigenvectors ( $\mathbf{A}\mathbf{v} = \mathbf{A}^T\mathbf{v} = \lambda\mathbf{v}$ )
- Why there is complex eigenvalue => because for rotation matrix, only times a real number, the direction of each vector are not changed

- Eigendecomposition

- Let  $\mathbf{V} \in \mathbb{R}^{n \times n}$  be a matrix whose column  $\mathbf{v}_i$  are  $n$  linearly independent eigenvectors of  $\mathbf{A}$  with  $\boldsymbol{\Lambda}$  the corresponding diagonal matrix of eigenvalues such that  $\boldsymbol{\Lambda}_{ii} = \lambda_i$ . Then:

$$\begin{aligned} \mathbf{AV} &= \sum_{j=0}^n \mathbf{Av}^{(j)} = \sum_{j=0}^n \lambda_j \mathbf{v}^{(j)} = \mathbf{V}\boldsymbol{\Lambda} \\ \mathbf{A} &= \mathbf{V}\boldsymbol{\Lambda}\mathbf{V}^{-1} \\ \mathbf{V}^{-1}\mathbf{AV} &= \boldsymbol{\Lambda} \end{aligned} \tag{20}$$

- Only *diagonalizable matrices* have **eigendecomposition** and **eigenvalues** (!!!eigenvalue is different from singular value)
- If  $\mathbf{A}$  is symmetric, we can choose  $n$  orthonormal eigenvectors so that  $\|\mathbf{v}_i\|_2 = 1$ ,  $\mathbf{v}_i^T \mathbf{v}_j = 0$  and  $n$  real eigenvalues  $\lambda_n \in \mathbb{R}$ . As a result, we have:

$$\begin{aligned} \mathbf{A} &= \mathbf{U}\boldsymbol{\Lambda}\mathbf{V}^T = \mathbf{A}^T = \mathbf{V}\boldsymbol{\Lambda}^T\mathbf{U}^T = \mathbf{V}\boldsymbol{\Lambda}\mathbf{U}^T (m = n) \\ \mathbf{U} &= \mathbf{V} \\ \mathbf{A} &= \mathbf{V}\boldsymbol{\Lambda}\mathbf{V}^{-1} = \sum_{i=1}^n \lambda_i \mathbf{v}_i \mathbf{v}_i^T \end{aligned} \tag{21}$$

- eigenvalues and singular values are the same for symmetric matrices
- If symmetric it is diagonalizable, if diagonalizable it may not be symmetric

- Representation of a vector in the Eigen Basis

- Similarly, if  $\mathbf{a}$  is an arbitrary vector, then we can also represent  $\mathbf{a}$  using the basis provided by the eigenvectors  $\mathbf{V}$  of a real symmetric matrix  $\mathbf{A}(\mathbf{a}\mathbf{a}^T)$ . We obtain:

$$\mathbf{a} = \sum_{i=1}^n \alpha_i \mathbf{v}_i \quad (22)$$

$$\mathbf{a}^T \mathbf{v}_i = \sum_{j \neq i} \alpha_j \mathbf{v}_j^T \mathbf{v}_i + \alpha_i \mathbf{v}_i^T \mathbf{v}_i = \alpha_i$$

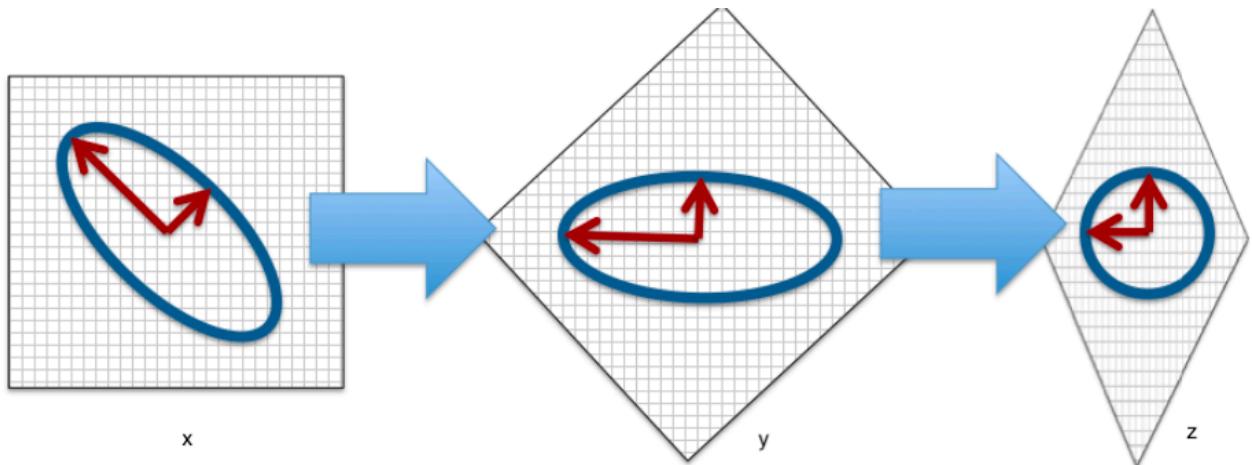
- Geometry: If  $A$  is a **real symmetric matrix** with *positive* eigenvalues, then the quadratic equation  $\mathbf{x}^T \mathbf{A} \mathbf{x} = 1$  defines an ellipsoid in an  $n$ -dimensional space, which provides a different way of thinking about these operations:

$$\mathbf{x}^T \mathbf{A} \mathbf{x} = 1, \mathbf{A} = \mathbf{V} \Lambda \mathbf{V}^T \quad (23)$$

$$\mathbf{x}^T \mathbf{V} \Lambda \mathbf{V}^T \mathbf{x} = 1, \text{ denotes } \mathbf{y} = \mathbf{V}^T \mathbf{x}$$

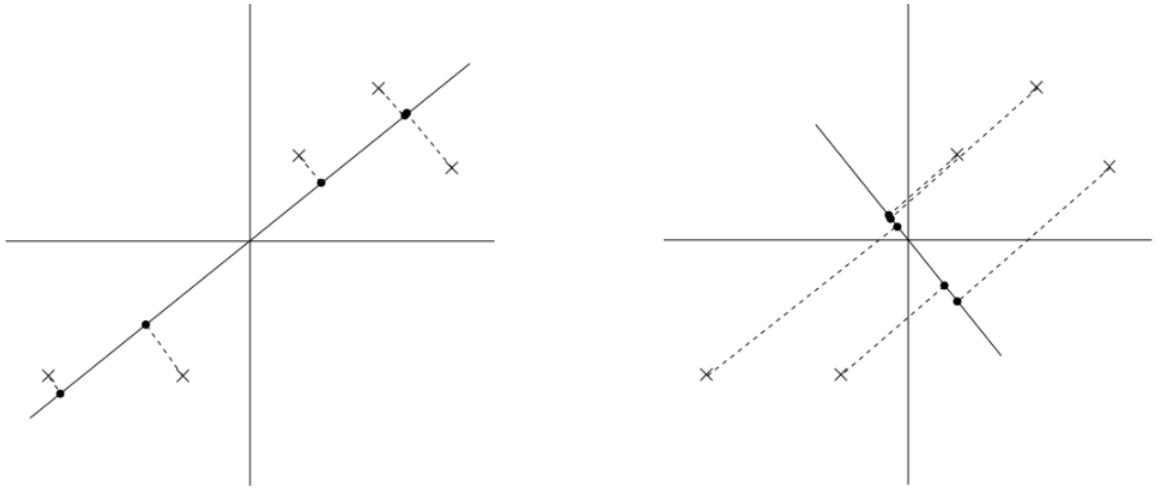
$$\mathbf{y}^T \Lambda \mathbf{y} = 1 \Rightarrow \mathbf{y}^T \Lambda^{\frac{1}{2}} \Lambda^{\frac{1}{2}} \mathbf{y} = 1, \text{ denotes } \mathbf{z} = \Lambda^{\frac{1}{2}} \mathbf{y}$$

$$\mathbf{z}^T \mathbf{z} = 1$$



- Principal Component Analysis (PCA) ([resource](#))

- Given a data matrix  $\mathbf{X} \in \mathbb{R}^{m \times n}$ , the goal of PCA is to identify the directions of maximum variance contained in the data (first one is better):
  - Project all the points to the direction s.t. the result data are well separated
  - Ignore the dimension where the data are similar (low variance), which means that this dimension is less useful



- Sample Variance in a Given Direction

  - Let  $\mathbf{v} \in \mathbb{R}^n$  s.t.  $\|\mathbf{v}\|_2 = \sqrt{\mathbf{v}^T \mathbf{v}} = 1$  (unit vector)

  - The sample estimate of the variance in the direction  $\mathbf{v}$  is given by the expression:

$$\frac{1}{m} \sum_{i=1}^m (\mathbf{v}^T \mathbf{x}^{(i)} - \mu)^2, \text{ where } \mu = \frac{1}{m} \sum_{i=1}^m \mathbf{v}^T \mathbf{x}^{(i)} \quad (24)$$

  - $\mathbf{v}^T \mathbf{x}^{(i)}$  represents the length of projection when project  $\mathbf{x}^{(i)}$  on  $\mathbf{v}$

- Pre-centering: Under the assumption that the data are pre-centered so that

$\frac{1}{m} \sum_{i=1}^m \mathbf{v}^T \mathbf{x}^{(i)} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}^{(i)} = 0$ , (向量和为零则每个方向上分量为零) this expression simplifies to:

$$\frac{1}{m} \sum_{i=1}^m (\mathbf{v}^T \mathbf{x}^{(i)})^2 = \frac{1}{m} (\mathbf{Xv})^T (\mathbf{Xv}) = \frac{1}{m} \mathbf{v}^T \mathbf{X}^T \mathbf{Xv} \quad (25)$$

- The Direction of the Maximum Variance

  - Suppose we want to identify the direction  $\mathbf{v}_1$  of maximum variance given the data matrix  $\mathbf{X}$ . We can formulate this optimization problem as follows (this is not a convex optimization => weak duality):

$$\begin{aligned} & \max_{\mathbf{v}} \frac{1}{m} \mathbf{v}^T \mathbf{X}^T \mathbf{Xv} \\ & \text{subject to } \|\mathbf{v}\|_2^2 = 1 \end{aligned} \quad (26)$$

  - Letting  $\Sigma = \frac{1}{m} \mathbf{X}^T \mathbf{X}$  which is a symmetric matrix, we form the Lagrangian

$$L(\mathbf{v}, \lambda) = \mathbf{v}^T \Sigma \mathbf{v} + \lambda(1 - \|\mathbf{v}\|_2^2) \quad (27)$$

  - It can be proved that the eigenvalue of  $\Sigma$  is non-negative by using SVD
  - Taking the derivative of  $L$  w.r.t  $\mathbf{v}$  and setting it to zero, we have:

$$\begin{aligned} \frac{\delta L(\mathbf{v}, \lambda)}{\delta \mathbf{v}} &= 2\Sigma \mathbf{v} - 2\lambda \mathbf{v} = 0 \\ \Sigma \mathbf{v} &= \lambda \mathbf{v} \end{aligned} \quad (28)$$

  - It means the optimal  $\mathbf{v}$  appears when  $\mathbf{v}$  is an eigenvector of  $\Sigma$

- As  $\mathbf{v} \neq 0$ ,  $\mathbf{v}$  must be an eigenvector of  $\Sigma$  with eigenvalue  $\lambda$ . Assuming  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  are the eigenvectors of  $\Sigma$ , corresponding to eigenvalues  $\sigma_1 \geq \dots \geq \sigma_b \geq 0$ , respectively, we choose the best  $\mathbf{v}$  while fixing  $\mathbf{X}$  we have (i.e. the largest mean value of  $\lambda_i$  where  $\lambda_i$  is the **singular values** of  $\mathbf{X}$ )

$$\begin{aligned}\mathbf{v}^* &= \mathbf{v}_1, \\ p^* &= \mathbf{v}_1^T \Sigma \mathbf{v}_1 = \lambda = \sigma_1\end{aligned}\tag{29}$$

- $l$  largest Directions of Variance

- Suppose instead of just the direction of maximum variance, we want the  $l$  largest directions of variance that are all mutually *orthogonal*
- Finding the second-largest direction of variance corresponds to solving the problem:

$$\begin{aligned}max_{\mathbf{v}} \mathbf{v}^T \Sigma \mathbf{v} \\ subject to \|\mathbf{v}\|_2^2 = 1 \\ \mathbf{v}^T \mathbf{v}_1 = 0\end{aligned}\tag{30}$$

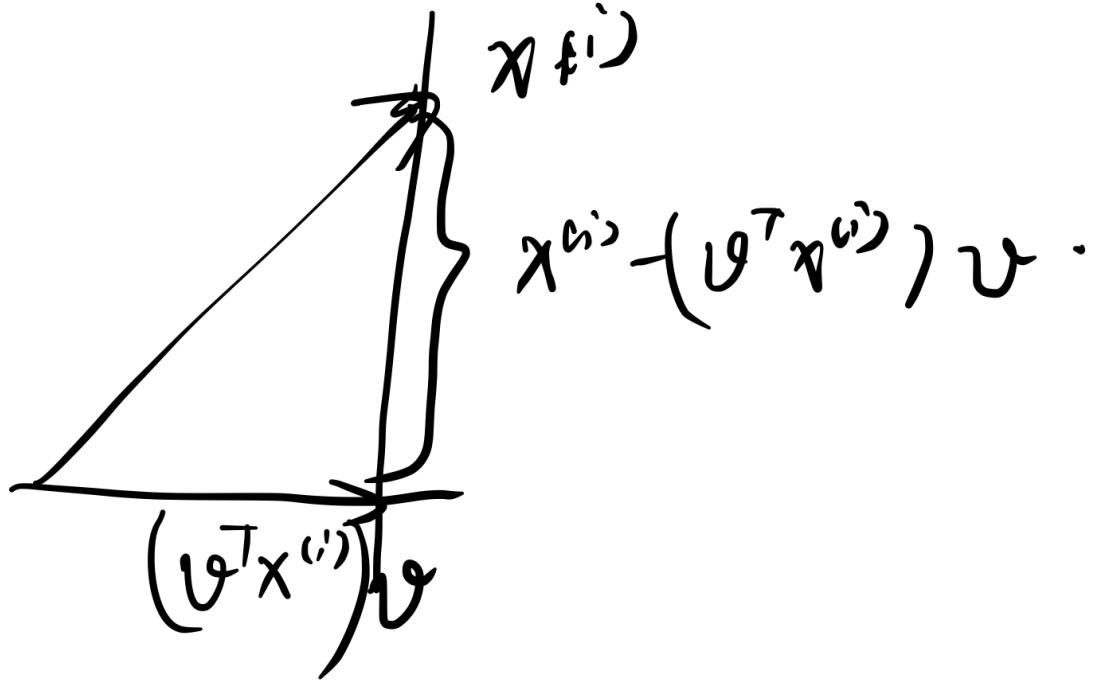
- We form the Lagrangian

$$\begin{aligned}L(\mathbf{v}, \lambda, v) &= \mathbf{v}^T \Sigma \mathbf{v} + \lambda(1 - \|\mathbf{v}\|_2^2) + v \mathbf{v}^T \mathbf{v}_1 \\ \frac{\delta L}{\delta \mathbf{v}} &= 2\Sigma \mathbf{v} - 2\lambda \mathbf{v} + v \mathbf{v}_1 = 0 \\ 2\mathbf{v}_1^T \Sigma \mathbf{v} - 2\lambda \mathbf{v}_1^T \mathbf{v} + v \mathbf{v}_1^T \mathbf{v}_1 &= 0 \\ 2(\Sigma(\Sigma = \Sigma^T) \mathbf{v}_1)^T \mathbf{v} - 2\lambda \mathbf{v}_1^T \mathbf{v} + v \mathbf{v}_1^T \mathbf{v}_1 &= 0 \\ 2\sigma_1 \mathbf{v}_1^T \mathbf{v} - 0 = v &= 0\end{aligned}\tag{31}$$

- In general, the top  $l$  directions of variance  $\mathbf{v}_1, \dots, \mathbf{v}_l$  are given by the  $l$  eigenvectors corresponding to the  $l$  largest eigenvalues of  $\frac{1}{m} \mathbf{X}^T \mathbf{X}$
- Minimize the Projection error:

- PCA can also be derived by picking the principal vectors that minimize the approximation error arising from projecting the data onto the  $l$ -dimensional subspace spanned by them

$$min_{\mathbf{v}} \frac{1}{m} \sum_{i=1}^m \|\mathbf{x}^{(i)} - (\mathbf{v}^T \mathbf{x}^{(i)}) \mathbf{v}\|_2^2\tag{32}$$



- Proof: Without loss of generality, we assume  $l = 1$ . The approximation error, i.e., the mean squared error, between  $\{\mathbf{x}(i)\}$  and their projections are

$$\begin{aligned}
& \frac{1}{m} \sum_{i=1}^m \|\mathbf{x}^{(i)} - (\mathbf{v}^T \mathbf{x}^{(i)}) \mathbf{v}\|_2^2 \\
&= \frac{1}{m} \sum_{i=1}^m (\|\mathbf{x}^{(i)}\|_2^2 + (\mathbf{v}^T \mathbf{x}^{(i)})^2 \|\mathbf{v}\|_2^2 - 2(\mathbf{v}^T \mathbf{x}^{(i)}) \mathbf{v}^T \mathbf{x}^{(i)}) \\
&= \frac{1}{m} \sum_{i=1}^m (\|\mathbf{x}^{(i)}\|_2^2 - (\mathbf{v}^T \mathbf{x}^{(i)})^2)
\end{aligned} \tag{33}$$

- Since  $\mathbf{x}_i$  is fixed for each  $i$ , minimizing the approximation error is equivalent to minimizing the variance of the projections:

$$\begin{aligned}
& \min_{\mathbf{v}} \frac{1}{m} \sum_{i=1}^m (-(\mathbf{v}^T \mathbf{x}^{(i)})^2) \\
&= \max_{\mathbf{v}} \frac{1}{m} \sum_{i=1}^m ((\mathbf{v}^T \mathbf{x}^{(i)})^2) \\
&= \max_{\mathbf{v}} \mathbf{v}^T \Sigma \mathbf{v}
\end{aligned} \tag{34}$$

- Dimensionality Reduction with PCA

- Data preprocessing: Compute  $\boldsymbol{\mu} = \frac{1}{m} \sum_i \mathbf{x}^{(i)}$  and replace each  $\mathbf{x}^{(i)}$  with  $\mathbf{x}^{(i)} - \boldsymbol{\mu}$ ; compute  $\sigma_j^2 = \frac{1}{m} \sum_i (x_j^{(i)})^2$  and replace each  $\mathbf{x}_j^{(i)}$  with  $\frac{\mathbf{x}_j^{(i)}}{\sigma_j}$  (normalize the data)
- Given pre-processed data matrix  $\mathbf{X} \in \mathbb{R}^{m \times n}$ , compute the sample covariance matrix  $\Sigma = \frac{1}{m} X^T X$
- Compute the  $l$  leading eigenvectors  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_l$  of  $\Sigma$  where  $\mathbf{v}_i \in \mathbb{R}^n$
- Stack the eigenvectors together into an  $n \times l$  matrix  $\mathbf{V}$  where each column  $i$  of  $\mathbf{V}$  corresponds

to  $\mathbf{v}_i$

- Project  $\mathbf{X}$  into the subspace spanned by the set of eigenvectors  $\mathbf{Z} = \mathbf{X}\mathbf{V}$
- To reconstruct  $\mathbf{X}$  given  $\mathbf{Z}$  and  $\mathbf{V}$ , we use  $\hat{\mathbf{X}} = \mathbf{Z}\mathbf{V}^T$

## Connection between PCA and SVD

- We have saw that the minimum Frobenius norm linear dimensionality reduction problem could be solved using the rank- $l$  SVD of  $\mathbf{X}$ :

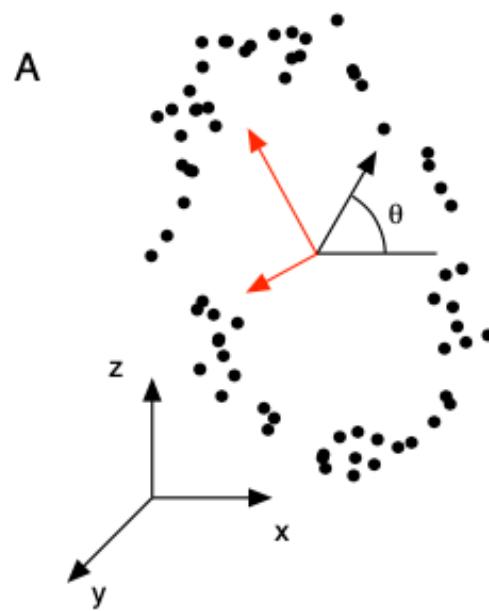
$$\operatorname{argmin}_{\mathbf{U}, \mathbf{S}, \mathbf{V}} \|\mathbf{X} - \mathbf{USV}^T\|_F^2 \quad (35)$$

- Where  $\mathbf{U} \in \mathbb{R}^{m \times l}$ ,  $\mathbf{S} \in \mathbb{R}^{l \times l}$ ,  $\mathbf{V} \in \mathbb{R}^{n \times l}$ . The matrix product  $\mathbf{Z} = \mathbf{US}$  gives the optimal rank- $l$  representation of  $\mathbf{X}$  with respect to Frobenius norm minimization
- Let  $l = n$  then  $\mathbf{X} = \mathbf{USV}^T$  and  $\mathbf{X}^T \mathbf{X} = \mathbf{VSU}^T \mathbf{USV}^T = \mathbf{VS}^2 \mathbf{V}^T$
- This means
  - right singular vectors of  $\mathbf{X}$  is the same as all the singular vectors of  $\mathbf{X}^T \mathbf{X}$
  - eigenvalues of  $\mathbf{X}^T \mathbf{X}$  are the diagonal elements of  $\mathbf{S}$
  - The  $l$  largest singular values for SVD and  $l$  largest eigenvalues for PCA correspond to the same  $l$  basis vectors

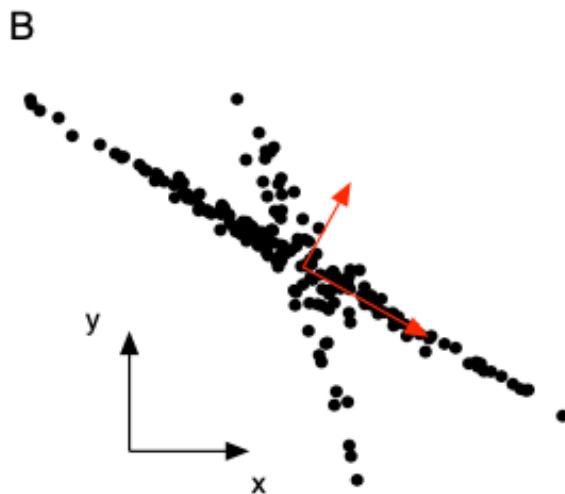
- According to PCA, the projection operation is  $\mathbf{Z} = \mathbf{X}\mathbf{V}$ , therefore

$$\mathbf{Z} = \mathbf{X}\mathbf{V} = (\mathbf{USV}^T)\mathbf{V} = \mathbf{US} \quad (36)$$

- Finally, note that if the decompositions are based only on the  $l$  leading principal vectors, which are identical under both PCA and SVD, the projection  $\mathbf{Z} = \mathbf{X}\mathbf{V}$  and  $\mathbf{Z} = \mathbf{US}$  will still be identical
- These manipulations show that PCA on  $\mathbf{X}^T \mathbf{X}$  and SVD on  $\mathbf{X}$  identify exactly the same sub space and result in exactly the same projection of the data into that sub-space
- As a result, generic linear dimensionality reduction simultaneously minimizes the Frobenius norm of the reconstruction error of  $\mathbf{X}$  and maximizes the retained variance in the learned sub-space
- Both SVD and PCA provide the same description of generic linear dimensionality reduction: an orthogonal basis for exactly the same optimal linear subspace
- Issues:
  - The computational complexity of PCA is  $O(n^2 m + n^3)$  if the full eigendecomposition is obtained and then truncated, compared to  $O(\min(nm^2, mn^2))$  for SVD
  - The basic SVD and PCA algorithm are not suitable for large-scale data. Instead, randomized algorithms are often used
  - The value of  $l$  can sometimes be chosen based on looking for eigenvalue gaps in the eigenspectrum of the covariance matrix
- When does PCA fail
  - The primary motivation behind PCA is to decorrelate the dataset, i.e., remove second-order dependencies. If higher-order dependencies exist between the features in the data, PCA may be insufficient at revealing all structure in the data (for example, the data distributed on a circle)



- PCA requires that each component must be perpendicular to the previous ones, but clearly this requirement is overly stringent and the data might be arranged along non-orthogonal axes



## Questions about Lagrange Multiplier

- When using weak duality, why it is ok to optimize L and f at the same time
- Why g is concave for lambda and v
- Why

## Lecture 10: Kernel Principal Component Analysis

- Limitation of PCA is it can only find the optimal linear subspace that retain maximum variances

## Basis Expansion

- One way to move beyond linear dimensionality reduction is to first apply a non-linear basis expansion to the data vectors, and then apply linear dimensionality reduction
- Any basis expansion can be applied (e.g. polynomials), just as in the case of classification and regression

- Given a data set  $\mathbf{X} \in \mathbb{R}^{m \times n}$  and a basis expansion function  $\phi : \mathbb{R}^n \mapsto \mathbb{R}^l$  for  $l > n$ , we obtain the following SVD-based algorithm:
  - Compute  $\mathbf{U}, \mathbf{S}, \mathbf{V} = SVD(\phi(\mathbf{X}))$
  - Return  $\mathbf{Z} = \mathbf{U}\mathbf{S}$

## Kernel PCA

- Step of Kernel PCA
  - Given a data set  $\mathbf{X} \in \mathbb{R}^{m \times n}$  and a basis expansion function  $\phi : \mathbb{R}^n \mapsto \mathbb{R}^l$  for  $l > n$ , we obtain the following PCA-based algorithm (expand to  $l$  and then reduce to  $k$ ):
    - Compute  $\Sigma = \frac{1}{m} \sum_i (\phi(\mathbf{x}^{(i)}) - \mu)(\phi(\mathbf{x}^{(i)}) - \mu)^T$  where  $\mu = \frac{1}{m} \sum_i \phi(\mathbf{x}^{(i)})$
    - Compute the  $k$  leading eigenvectors  $\mathbf{v}_1, \dots, \mathbf{v}_k$  of  $\Sigma$  where  $\mathbf{v}_j \in \mathbb{R}^{l \times 1}$  for  $j = 1, \dots, k$
    - Stack the eigenvectors together to form  $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k]$ , where  $\mathbf{V} \in \mathbb{R}^{l \times k}$
    - Project the matrix  $\phi(\mathbf{X})$  into the rank- $k$  subspace of maximum variance by computing the matrix product  $\mathbf{Z} = \phi(\mathbf{X})\mathbf{V}$
- Kernel function
  - As in the classification case, it becomes very expensive to use an explicit basis function expansion that maps data into a high-dimensional space
  - In the basic SVD-based algorithm, there's no way to avoid this problem
  - In the PCA-based algorithm, we are able to take advantage of the kernel trick, which provides a much more efficient computation of feature-based inner products without explicitly computing the basis expansion:

$$K(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \phi(\mathbf{x}^{(i)})^T \phi(\mathbf{x}^{(j)}) \quad (37)$$

- Convert the inner product used in kernel PCA

- Given  $\phi : \mathbb{R}^n \mapsto \mathbb{R}^l$ , we compute the covariance matrix in the new feature space

$$\Sigma = \frac{1}{m} \sum_{j=1}^m \phi(\mathbf{x}^{(j)}) \phi(\mathbf{x}^{(j)})^T \quad (38)$$

- Eigendecomposition of  $\Sigma$  is given by (since it is a symmetric matrix):

$$\Sigma \mathbf{v}_k = \frac{1}{m} \sum_{j=1}^m \phi(\mathbf{x}^{(j)}) \phi(\mathbf{x}^{(j)})^T \mathbf{v}_k = \lambda_k \mathbf{v}_k, \forall k = 1, \dots, l \quad (39)$$

- It is not hard to see that  $\mathbf{v}_k$  can be expressed as

$$\begin{aligned}
\mathbf{v}_k &= \frac{1}{m\lambda_k} \sum_{j=1}^m \phi(\mathbf{x}^{(j)}) \phi(\mathbf{x}^{(j)})^T \mathbf{v}_k \\
&= \sum_{j=1}^m \frac{1}{m\lambda_k} [\phi(\mathbf{x}^{(j)})^T \mathbf{v}_k] \phi(\mathbf{x}^{(j)}) (\text{scalar}) \\
&= \sum_{j=1}^m w_k^{(j)} \phi(\mathbf{x}^{(j)}) \text{ where } w_k^{(j)} = \frac{1}{m\lambda_k} [\phi(\mathbf{x}^{(j)})^T \mathbf{v}_k] \\
\phi(\mathbf{x}^{(i)})^T \mathbf{v}_k &= \sum_{j=1}^m w_k^{(j)} \phi(\mathbf{x}^{(i)})^T \phi(\mathbf{x}^{(j)}) = m\lambda_k w_k^{(j)} \\
\sum_{j=1}^m w_k^{(j)} K_{ij} &= m\lambda_k w_k^{(j)} \\
\mathbf{K}\mathbf{w}_k &= m\lambda_k \mathbf{w}_k
\end{aligned} \tag{40}$$

- We can find that  $\mathbf{w}_k$  is the  $k$ th eigenvector of  $\mathbf{K}$  with the corresponding eigenvalue  $m\lambda_k$
- As the eigenvalues of  $\mathbf{K}$  are proportional to the eigenvalues  $\lambda'$ s of the covariance matrix  $\Sigma$  in the feature space, PCA on  $\Sigma \in \mathbb{R}^{l \times l}$  is equivalent to PCA on  $\mathbf{K} \in \mathbb{R}^{m \times m}$
- For any new data point  $\mathbf{x}$  with its high-dimensional feature representation  $\phi(\mathbf{x})$ , the  $k$ th PCA component can be obtained by projecting it on the  $k$ -th eigenvector  $\mathbf{v}_k$  of  $\Sigma$

$$\phi(\mathbf{x})^T \mathbf{v}_k = \sum_{i=1}^m w_k^{(i)} \phi(\mathbf{x})^T \phi(\mathbf{x}^{(i)}) = \sum_{j=1}^m w_k^{(i)} K(\mathbf{x}, \mathbf{x}_i) \tag{41}$$

- Complete steps of PCA: Given a data set  $\mathbf{X} \in \mathbb{R}^{m \times n}$  and a kernel function  $K$ , kernel PCA can be computed as follows:
  - Compute  $K_{ij} = \mathbf{x}_i, \mathbf{x}_j$  for all  $i, j$
  - Compute  $\mathbf{K}' = (\mathbf{I} - \mathbf{1}_m)\mathbf{K}(\mathbf{I} - \mathbf{1}_m)$  where  $\mathbf{1}_m$  is an  $m \times m$  matrix where every entry is  $\frac{1}{m}$  (to zero centro the data)

$$\begin{aligned}
\phi'(\mathbf{x}^{(i)}) &= \phi(\mathbf{x}^{(i)}) - \frac{1}{m} \sum_{k=1}^m \phi(\mathbf{x}^{(k)}) \\
\mathbf{K}'_{ij} &= \phi'(\mathbf{x}^{(i)})^T \phi'(\mathbf{x}^{(j)}) \\
&= [\phi(\mathbf{x}^{(i)}) - \frac{1}{m} \sum_{k=1}^m \phi(\mathbf{x}^{(k)})]^T [\phi(\mathbf{x}^{(j)}) - \frac{1}{m} \sum_{k=1}^m \phi(\mathbf{x}^{(k)})] \\
&= \mathbf{K}_{ij} - \frac{1}{m} \sum_{k=1}^m \mathbf{K}_{ik} - \frac{1}{m} \sum_{k=1}^m \mathbf{K}_{kj} + \frac{1}{m^2} \sum_{k=1}^m \sum_{k'=1}^m \mathbf{K}_{kk'} \\
&= \mathbf{K}_{ij} - \mathbf{K}_{i,*} \mathbf{1}_m - \mathbf{1}_m^T \mathbf{K}_{*,j} + \mathbf{1}_m^T \mathbf{K} \mathbf{1}_m \\
\mathbf{K} &= \mathbf{K} - \mathbf{K} \mathbf{1}_{m \times m} - \mathbf{1}_{m \times m} \mathbf{K} + \mathbf{1}_{m \times m}^T \mathbf{K} \mathbf{1}_{m \times m} \\
&= (\mathbf{I} - \mathbf{1}_m) \mathbf{K} (\mathbf{I} - \mathbf{1}_m)
\end{aligned} \tag{42}$$

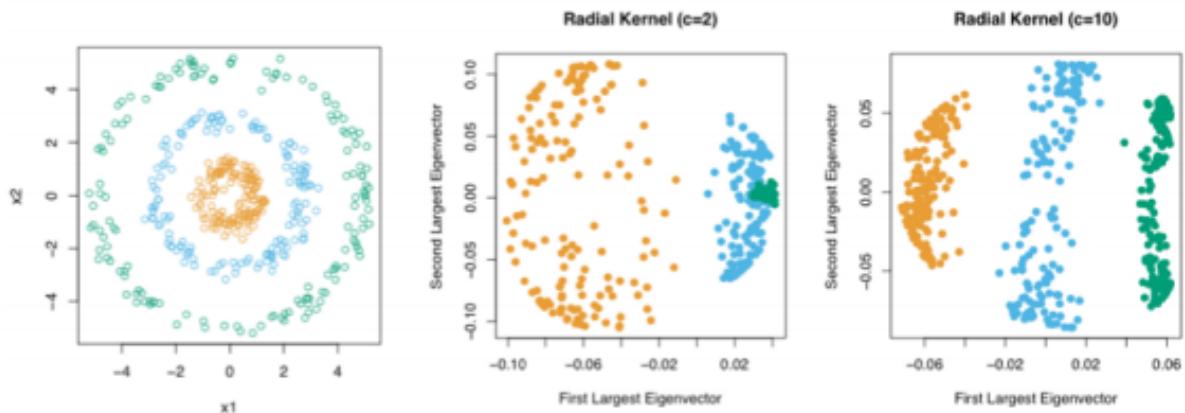
- Compute the  $r$  leading eigenvectors  $\mathbf{w}_1, \dots, \mathbf{w}_r$  of  $\mathbf{K}'$  along with their eigenvalue  $m\lambda_1, \dots, m\lambda_r$

- Compute the  $k$ th element of the projected data vector  $\mathbf{z} \in \mathbb{R}^{r \times 1}$  using  

$$z_k = \sum_{j=1}^m w_k^{(i)} K(\mathbf{x}, \mathbf{x}_i)$$
- Example:

**Consider the case of using the radial basis function kernel**

$$\mathcal{K}(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{1}{c} \|\mathbf{x} - \mathbf{y}\|_2^2\right)$$



- Summary

- Kernel PCA provides a non-linear dimensionality reduction method that can be used to extract directions of maximum variance without explicitly performing basis expansion
- Kernel PCA can be thought of as looking for directions of variance in the space of similarities between data cases, and can extract components that correspond to fairly complex structures
- When a good kernel is known a priori, kernel PCA can provide an effective pre-processing step for clustering methods as well as linear classification and regression methods
- However, exact computation of kernel PCA can be expensive because the size of the matrix to be decomposed is  $m \times m$

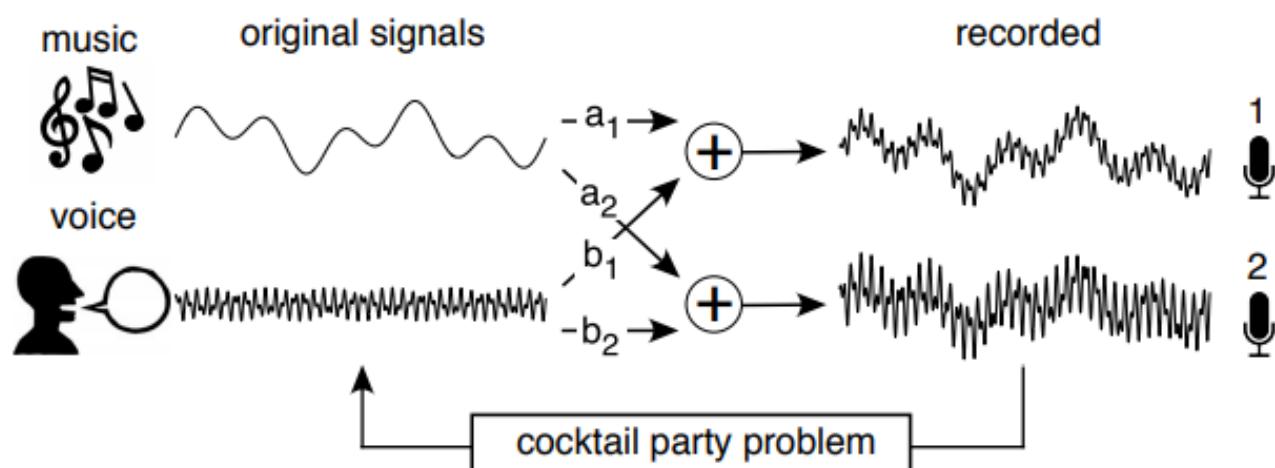
# Lecture 11: Independent Component Analysis

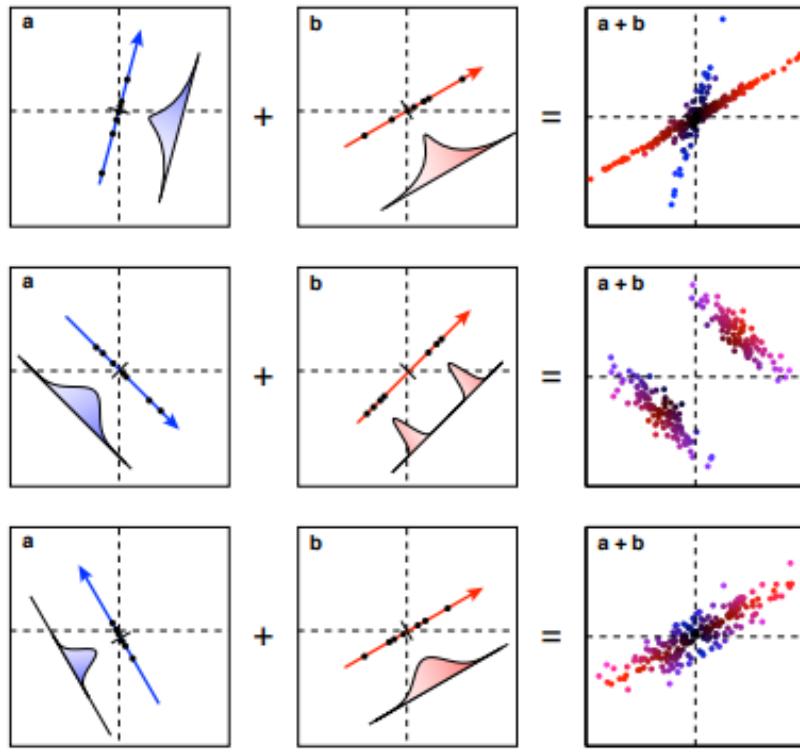
## Independent Component Analysis Motivation

- Measurements in the real world are corrupted by random noise - but that is only one piece of the story. Often, measurements cannot be made in isolation, but reflect the combination of many distinct sources
- Linear mixing of independent sources is exactly what occurs when you listen to multiple audio sources at the same time
- Humans are able to de-mix multiple sources of audio (multiple people speaking) into distinct sources very accurately
- Independent component analysis (ICA) is designed to solve exactly this problem (called blind source separation) and can do so very reliably when the number of observed linearly mixed channels is equal to the number of sources
- [Demo](#)

## ICA Illustration

- An example





- Where the  $x$ -axis and the  $y$ -axis represents the signal of music recorded for microphone 1 and 2 respectively (the loud of the sound times the corresponding weight for two microphones)
- Since the voice has positive correlation to the distance between the source and the microphone, the proportion of voices recorded by different microphone is a constant because of the fixed distance => if each data point represents the record of a specific time, the points for the same source are distributed on a line throw the origin with the slope equals to the proportion of distances between source and two different microphones
- We assume that the points for each microphone has a normal distribution on time (the signal is time independent)
- In this case, it is useless to find a maximum variance direction

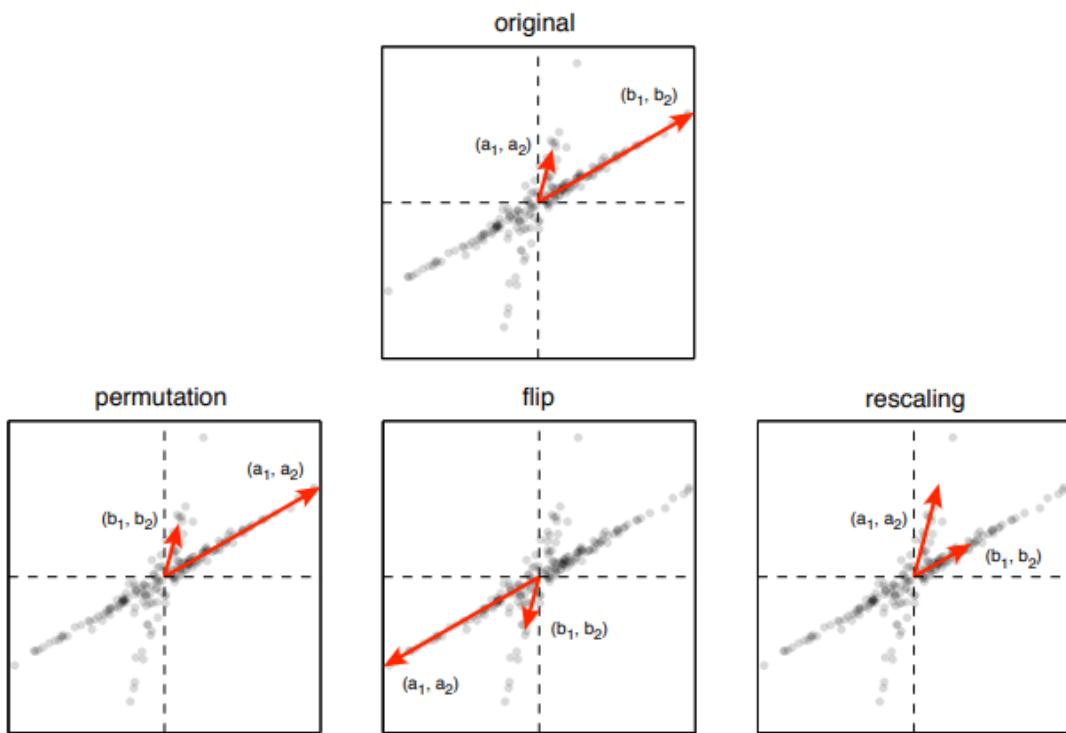
## ICA Formulation

- To formalize this problem, assume the data  $\mathbf{s} \in \mathbb{R}^{n \times 1}$  is generated via  $n$  independent sources (at a specific time). What we observe is
$$\mathbf{x} = \mathbf{As} \quad (1)$$
  - where  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is an unknown *mixing matrix*
  - **Each row of  $\mathbf{A}$  represents a microphone**
- Repeated observations (get the sample different time) gives us a dataset  $D : \{\mathbf{x}^{(i)}, i = 1, \dots, m\}$ , and our goal is to find  $\mathbf{A}$ , so that given the observation  $\mathbf{x}^{(i)}$ , we can recover the sources  $\mathbf{s}^{(i)}$  by computing  $\mathbf{s}^{(i)} = \mathbf{Wx}^{(i)}$ , where  $\mathbf{W} = \mathbf{A}^{-1}$  is the *unmixing matrix*
- In our cocktail party problem,  $s_j^{(i)}$  is the sound that speaker  $j$  was uttering at time  $i$ . Similarly,  $\mathbf{x}_j^{(i)}$  is acoustic reading recorded by microphone  $j$  at time  $i$
- Notice that here the  $i$  is the index of **each column**

## ICA Ambiguities (歧義, 等价)

- To what degree can  $\mathbf{W}^{-1} = \mathbf{A}$  be recovered
- Permutation ambiguity: The labels of each independent component can be arbitrarily permuted. That is, given only the  $\mathbf{x}^{(i)}$ 's, there will be no way to distinguish between  $\mathbf{A}$  and  $\mathbf{AP}$ , where  $\mathbf{P} \in \mathbb{R}^{n \times n}$  is a permutation matrix (the order of independent components can be changed)
- Sign ambiguity: Any independent component can be flipped across the origin. That is,  $\mathbf{A}$  and  $-\mathbf{A}$  are treated identical in practical applications (the direction of independent components can be changed)
- Scale ambiguity: Any independent component can be rescaled with arbitrary length. For example, given  $\alpha \neq 0$ , we have (the length of independent components can be changed, and absorbed to the matrix  $\mathbf{A}$ )

$$\mathbf{x}^{(i)} = \mathbf{As}^{(i)} = \alpha \mathbf{A} \left( \frac{s^{(i)}}{\alpha} \right) = \hat{\mathbf{A}} \hat{\mathbf{s}}^{(i)} \quad (2)$$



## A Strategy for Solving ICA

- Divide-and-conquer provides a strategy to solve this problem. Rather than trying to solve for  $\mathbf{s}$  and  $\mathbf{A}$  simultaneously, we focus on finding  $\mathbf{A}$ . Furthermore, rather than trying to solve for  $\mathbf{A}$  all at once, we solve for  $\mathbf{A}$  in a piece-meal fashion by cutting up  $\mathbf{A}$  into simpler and more manageable parts
- SVD of  $\mathbf{A}$  is given by  $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$  (unknown)
- We estimate  $\mathbf{A}$  and its inverse  $\mathbf{W}$  by recovering each piece of the decomposition individually

$$\mathbf{W} = \mathbf{A}^{-1} = \mathbf{V}\Sigma^{-1}\mathbf{U}^T \quad (3)$$

- To compute  $\mathbf{U}$  and  $\Sigma$ , we first estimate the covariance of the observed mixed data in terms of the underlying sources (PCA)

$$\begin{aligned}
\mathbf{C} &= \frac{1}{m} \sum_{i=1}^m \mathbf{x}^{(i)} (\mathbf{x}^{(i)})^T = \frac{1}{m} \sum_{i=1}^m \mathbf{A} \mathbf{s}^{(i)} (\mathbf{A} \mathbf{s}^{(i)})^T \\
&= \frac{1}{m} \sum_{i=1}^m \mathbf{U} \Sigma \mathbf{V}^T \mathbf{s}^{(i)} (\mathbf{U} \Sigma \mathbf{V}^T \mathbf{s}^{(i)})^T \\
&= \mathbf{U} \Sigma \mathbf{V}^T \left( \frac{1}{m} \sum_{i=1}^m \mathbf{s}^{(i)} (\mathbf{s}^{(i)})^T \right) \mathbf{V} \Sigma \mathbf{U}^T \\
&= \mathbf{U} \Sigma^2 \mathbf{U}^T,
\end{aligned}$$

assuming  $\frac{1}{m} \sum_i \mathbf{x}^{(i)} = 0$  and  $\frac{1}{m} \sum_i \mathbf{s}^{(i)} (\mathbf{s}^{(i)})^T = \mathbf{I}$

- The first assumption get the equation and the second one is from the three ambiguities

$$\begin{aligned}
\frac{1}{m} \sum_i \mathbf{s}^{(i)} \mathbf{s}^{(i)T} &= E(\mathbf{s}^{(i)} \mathbf{s}^{(i)T}) = \\
\begin{bmatrix} E(\mathbf{s}_1 \mathbf{s}_1) & E(\mathbf{s}_1 \mathbf{s}_2) & \dots & E(\mathbf{s}_1 \mathbf{s}_n) \\ \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & E(\mathbf{s}_n \mathbf{s}_n) \end{bmatrix} &= \mathbf{I} = \frac{1}{m} \mathbf{s} \mathbf{s}^T
\end{aligned} \tag{4}$$

- Each row is orthogonal to other rows
- Because we assume  $\mathbf{s}$  is normalized and different features (**speakers**) of  $\mathbf{s}$  are independent
  - $E(s_i s_j) = 0$  when  $i \neq j$

$$\begin{aligned}
E(S_i S_j) &= \sum_{s_i \in S_i, s_j \in S_j} p(S_i = s_i, S_j = s_j) s_i s_j \\
&= \sum_{s_i \in S_i, s_j \in S_j} p(S_i = s_i) p(S_j = s_j) s_i s_j \\
&= \sum_{s_i \in S_i} p(S_i = s_i) s_i \sum_{s_j \in S_j} p(S_j = s_j) s_j \\
&= E(S_i) E(S_j) = 0 * 0 = 0
\end{aligned} \tag{5}$$

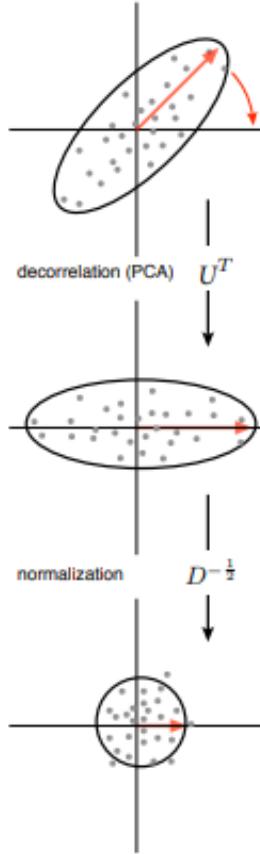
- $E(s_i s_j) = 1$  when  $i = j$

$$\begin{aligned}
Var(S_i) &= \sum_{s_i \in S_i} (s_i - \mu)^2 p(s_i) \\
&= \sum_{s_i \in S_i} (s_i^2 + \mu^2 - 2s_i\mu)p(s_i) \\
&= \sum_{s_i \in S_i} s_i^2 p(s_i) + \sum_{s_i \in S_i} \mu^2 p(s_i) - 2 \sum_{s_i \in S_i} s_i \mu p(s_i) \\
&= \sum_{s_i \in S_i} s_i^2 p(s_i) + \mu^2 \sum_{s_i \in S_i} p(s_i) - 2\mu \sum_{s_i \in S_i} s_i p(s_i) \\
&= \sum_{s_i \in S_i} s_i^2 p(s_i) + \mu^2 - 2\mu^2 \\
&= E(S_i^2) - E(S_i)^2 \\
Var(S_i) &= 1, E(S_i) = 0 \\
E(S_i^2) &= 1
\end{aligned} \tag{6}$$

- Transform the problem
  - By our shrewd(精明) choice of assumption, the covariance of the data  $\mathbf{C} = \mathbf{U}\Sigma^2\mathbf{U}^T$  is independent of sources  $\mathbf{s}^{(i)}$ 's as well as  $\mathbf{V}$ , which can be solved using eigendecomposition
  - Therefore, if our assumptions behind ICA are correct, then we have identified a partial solution to **A**:  $\mathbf{U}$  is a matrix of the **stacked eigenvectors of the covariance of the data ( $\mathbf{x}$ )** and  $\Sigma$  is a diagonal matrix with the square root of the associated eigenvalue in the diagonal
  - Interestingly, the eigenvectors of the covariance are essentially the principal components of the data in PCA
- Defining  $\hat{\mathbf{x}}^{(i)} = (\Sigma^{-1}\mathbf{U}^T)\mathbf{x}^{(i)} = \mathbf{V}^T\mathbf{s}^{(i)}$ , we have:

$$\hat{\mathbf{C}} = \frac{1}{m} \sum_{i=1}^m \hat{\mathbf{x}}^{(i)} (\hat{\mathbf{x}}^{(i)})^T = \frac{1}{m} \sum_{i=1}^m \mathbf{V}^T \mathbf{s}^{(i)} (\mathbf{V}^T \mathbf{s}^{(i)})^T = \mathbf{I} \tag{7}$$

- where  $\mathbf{x}^{(i)}$  is first projected on the principal components,  $\mathbf{U}^T\mathbf{x}^{(i)}$ , to remove linear correlations, and then scaled so that every direction has unit variance  $\Sigma^{-1}\mathbf{U}^T\mathbf{x}^{(i)}$ . In signal processing, this is termed as *whitening*



- Plugging  $\hat{\mathbf{x}}^{(i)}$  into  $s^{(i)} = \mathbf{W}\mathbf{x}^{(i)}$ , we simplify ICA down to solving  $\mathbf{s}^{(i)} = \mathbf{V}\hat{\mathbf{x}}^{(i)}$ , i.e. finding a rotation matrix  $\mathbf{V}$
  - We now need to **exploit the statistics of independence** to identify  $\mathbf{V}$ . Whitening the data removed all second-order correlations, hence discerning the last rotation matrix  $\mathbf{V}$  requires examining other measures of dependency
  - Statistical independence is the strongest measure of dependency between random variables. In particular, if two random variables  $a$  and  $b$  are independent, then  $p(a, b) = p(a)p(b)$  - the joint probability factorizes. In the context of ICA, we assume that all sources are statistically independent, thus  $p(\mathbf{s}) = \prod_j p(\mathbf{s}_j)$
  - Mutual Information
    - We resort to *multi-information*, a generalization of mutual information from information theory, to **judge how close a distribution is to statistical independence for multiple variables**
- $$I(\mathbf{s}) = \sum_{s \in S} p(s) \log \frac{p(s)}{\prod_j p(s_j)} \quad (8)$$
- It is a non-negative quantity that reaches a minimum of zero if and only if all variables are statistically independent. For example, if  $p(\mathbf{s}) = \prod_j p(\mathbf{s}_j)$ ,  $I(\mathbf{s}) = 0$
  - Minimizing the multi-information is difficult in practice but can be simplified
  - The multi-information is a function of the entropy, which is defined as  $H(\mathbf{s}) = -\sum_{s \in S} p(s) \log p(s)$ , measuring the amount of uncertainty about a distribution  $p(\mathbf{s})$
  - The multi-information is the difference between the sum of entropies of the marginal distributions and the entropy of the joint distribution
  - $|\mathbf{V}| = 1$ , since  $\mathbf{V}$  is orthogonal

$$\begin{aligned}
I(\mathbf{s}) &= \sum_j H(s_j) - H(\mathbf{s}) = \sum_j H((\mathbf{V}\hat{\mathbf{x}})_j) - H(\mathbf{V}\hat{\mathbf{x}}) \\
&= \sum_j H((\mathbf{V}\hat{\mathbf{x}})_j) - (H(\hat{\mathbf{x}}) + \log_2 |\mathbf{V}|) = \sum_j H((\mathbf{V}\hat{\mathbf{x}})_j) - H(\hat{\mathbf{x}})
\end{aligned} \tag{9}$$

### 9.3 Entropy under Linear Transformations (Optional)

For a linear transformation  $A \in \mathbb{R}^{n \times n}$  and a random variable  $x \in \mathbb{R}^{n \times 1}$ , we have

$$H[Ax] = H[x] + \log_2 |A| \tag{173}$$

where  $|A|$  is the absolute value of the determinant<sup>25</sup> of  $A$ , and is equal to the absolute value of the product of its eigenvalues.

*Proof.* Let  $z = Ax$ . The input and output probabilities are related by

$$p(z)dz = p(x)dx, \tag{174}$$

on an infinitesimal area and

$$p(z) = \frac{p(x)}{\left| \frac{dz}{dx} \right|} = \frac{p(x)}{|A|}. \tag{175}$$

The entropy of  $z$  can be computed by

$$\begin{aligned}
H[z] &= - \int p(z) \log_2 p(z) dz \\
&= - \int \frac{p(x)}{|A|} \log_2 \frac{p(x)}{|A|} |A| dx \\
&= - \int p(x) \log_2 \frac{p(x)}{|A|} dx \\
&= - \int p(x) \log_2 p(x) dx + \int p(x) \log_2 |A| dx \\
&= H[x] + \log_2 |A|.
\end{aligned} \tag{176}$$

□

- Optimization

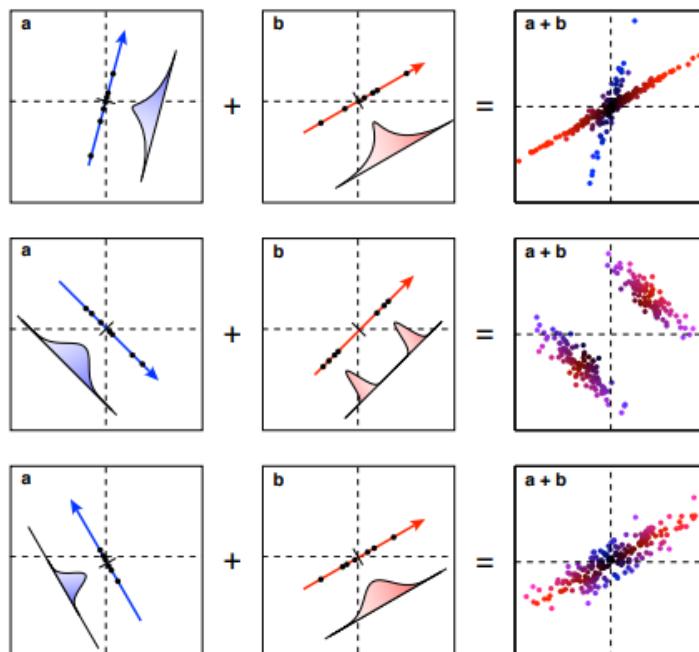
- The optimization is simplified further by recognizing that the term  $H(\hat{\mathbf{x}})$  is a constant, independent of  $V$ , and can be dropped:

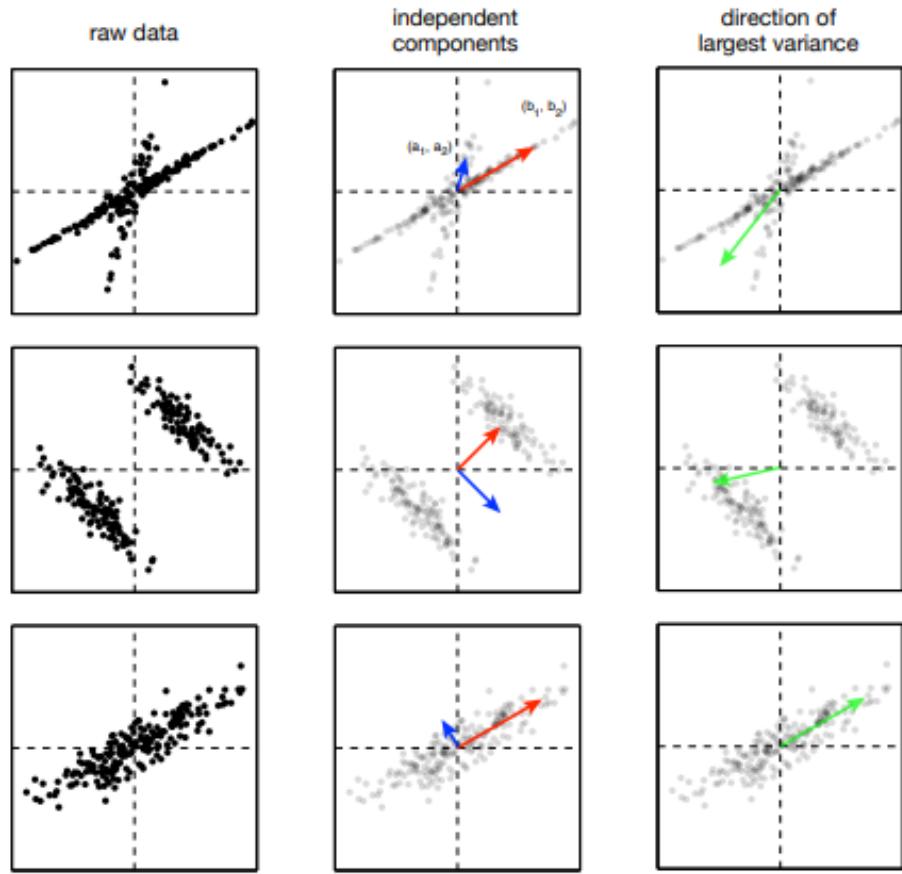
$$\mathbf{V}^* = \arg \min_{\mathbf{V}} \sum_j H((\mathbf{V}\hat{\mathbf{x}})_j)$$

- Calculating the entropy from a finite data set is difficult and should be approached with abundant caution. **Many ICA variants focus on approximations to the above equation**
- It is equivalent to finding the rotation that maximizes the expected log-likelihood of the observed data under the assumption that the data are statistically independent

- Likewise, a solution to the equation finds the rotation that maximizes the “non-Gaussianity” of the transformed data
- The ICA Algorithm
  - Subtract off the mean of the data in each dimension
  - Whiten the data by calculating the **eigenvectors of the covariance of the data**
  - Identify final rotation matrix that optimizes statistical independence

## PCA vs ICA





- Discussion
  - By assumption ICA expects the data to arise from an **(invertible) linear transformation** applied to a **factorial distribution**. Hence, ICA expects and will work optimally on data arising from linearly transformed factorial distributions - applying ICA to any other type of data is a gamble with no guarantees
  - ICA is equivalent to PCA for Gaussian data because a whitening filter based on PCA already achieves a factorial distribution and no need exists for computing  $V$
  - If higher-order correlations are small or insignificant, then PCA is sufficient to recover the independent sources of the data
  - Empirically measuring  $I(s)$  provides a useful check on how independent the recovered sources are. In practice,  $I(s)$  rarely achieves 0 either because of sampling issues or because of local minima in the optimization

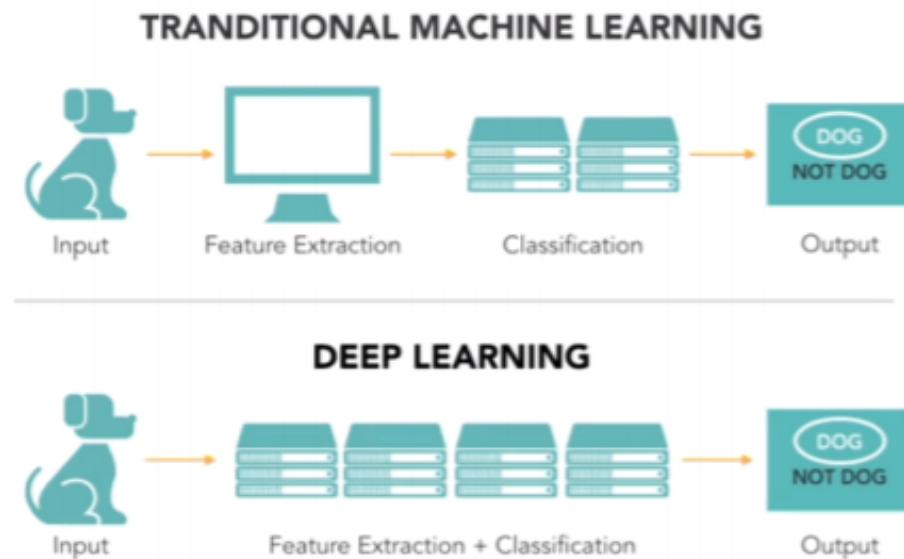
## Questions

- What is the second order correlation
- Why PCA and ICA are equivalents for Gaussian data

# Lecture 12: Neural Networks and Deep Learning

## Overview

- Feature Learning

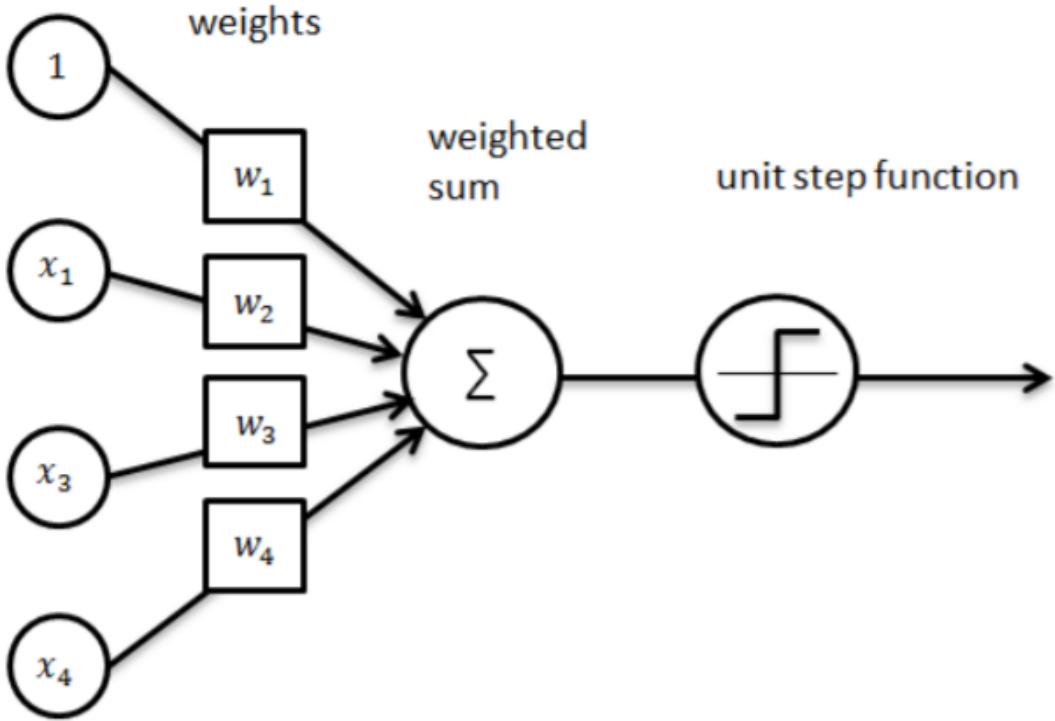


- Grand Challenge: Why/how deep learning works
  - Gap between engineering (or art) and science: Lack of theoretical understandings guarantees, and analytical tools
  - Training is computationally expensive and difficult, relying on many “magics”
    - Structure design
    - Network initialization
    - Hyper-parameter tuning
  - No principled way to incorporate(包括) domain expertise(专长)
  - No principled way to interpret the model behaviors

## Neural Networks

- McCulloch and Pitts Neuron (1943)

inputs



- The Perceptron (感受器) (for **binary** classification problems)

- The Perceptron is a simple online algorithm for adapting the weights in a McCulloch/Pitts neuron. It was developed in the 1950s by Rosenblatt at Cornell
  - Start with all-zero weight vector  $\mathbf{w}^{(0)}$ , and initialize  $t$  to 0
  - For all  $(\mathbf{x}^{(i)}, y^{(i)}) \in D$ , compute the activation  $\mathbf{a}^{(i)} = (\mathbf{x}^{(i)})^T \mathbf{w}^{(t)}$
  - If  $y^{(i)} \mathbf{a}^{(i)} < 0$ , then  $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + y^{(i)} \mathbf{x}^{(i)}$  and  $t \leftarrow t + 1$

$$\begin{aligned}\mathbf{a}_{t+1}^{(i)} &= (\mathbf{x}^{(i)})^T \mathbf{w}^{(t+1)} \\ &= (\mathbf{x}^{(i)})^T (\mathbf{w}^{(t)} + y^{(i)} \mathbf{x}^{(i)}) \\ &= a_t^{(i)} + y^{(i)} (\mathbf{x}^{(i)})^T \mathbf{x}^{(i)}\end{aligned}\tag{1}$$

$$\hat{y} = \begin{cases} 1, & \text{if } \mathbf{w}^T \mathbf{x} \geq 0 \\ -1 & \text{if } \mathbf{w}^T \mathbf{x} < 0, \end{cases}$$

- If  $y^{(i)} = 1$  and  $\hat{y}^{(i)} = -1$ , the activation ( $\hat{y}^{(i)}$ ) is initially negative and will be increased (see the equation)
- If  $y^{(i)} = -1$  and  $\hat{y}^{(i)} = 1$ , the activation is initially positive and will be decreased

- The Perceptron Theorem

- Let the training set be  $\mathbf{D} = \{(\mathbf{x}^{(i)}, y^{(i)}), i = 1, \dots, m\}$ . Assume that  $\|\mathbf{x}^{(i)}\|_2 \leq \alpha$  for all  $i$  and further that there exists a unit-length vector  $\mathbf{u}$  ( $\|\mathbf{u}\|_2 = 1$ ) such that  $y^{(i)} (\mathbf{u}^T \mathbf{x}^{(i)}) \geq \beta$  (weighted average of dimensions  $\mathbf{x}^{(i)}$  times  $y^{(i)}$ ) for all examples in  $D$ . Then the total number of mistakes that the perceptron algorithm makes on this sequence is at most  $(\alpha/\beta)^2$

- o Proof:

- The perceptron updates its weights only for following case:

$$(x^{(i)})^T w^{(t)} y^{(i)} \leq 0 \quad (2)$$

- o According the rule of perceptron learning

$$u^T w^{(t+1)} = u^T (w^{(t)} + y^{(i)} x^{(i)}) \geq u^T w^{(t)} + \beta \quad (3)$$

- The weight matrix is initialized as 0 :

$$u^T w^{(t+1)} \geq u^T w^{(1)} + t\beta \geq t\beta \quad (4)$$

- o Therefore:

$$\begin{aligned} \|w^{(t+1)}\|_2^2 &= \|w^{(t)} + y^{(i)} x^{(i)}\|_2^2 \\ &= \|w^{(t)}\|_2^2 + \|y^{(i)} x^{(i)}\|_2^2 + 2y^{(i)} (x^{(i)})^T w^{(t)} \\ &= \|w^{(t)}\|_2^2 + \|x^{(i)}\|_2^2 + 2y^{(i)} (x^{(i)})^T w^{(t)} \\ &\leq \|w^{(t)}\|_2^2 + \|x^{(i)}\|_2^2 \\ &\leq \|w^{(t)}\|_2^2 + \alpha^2 \leq \|w^{(1)}\|_2^2 + t\alpha^2 = t\alpha^2 \end{aligned} \quad (5)$$

- o Therefore:

$$\begin{aligned} \sqrt{t}\alpha &= \|w^{(t+1)}\|_2 \\ &= \|w^{(t+1)}\|_2 \|u\|_2 \geq \|w^{(t+1)}\|_2 \|u\|_2 \cos \theta \\ &= u^T w^{(t+1)} \geq t\beta \end{aligned} \quad (6)$$

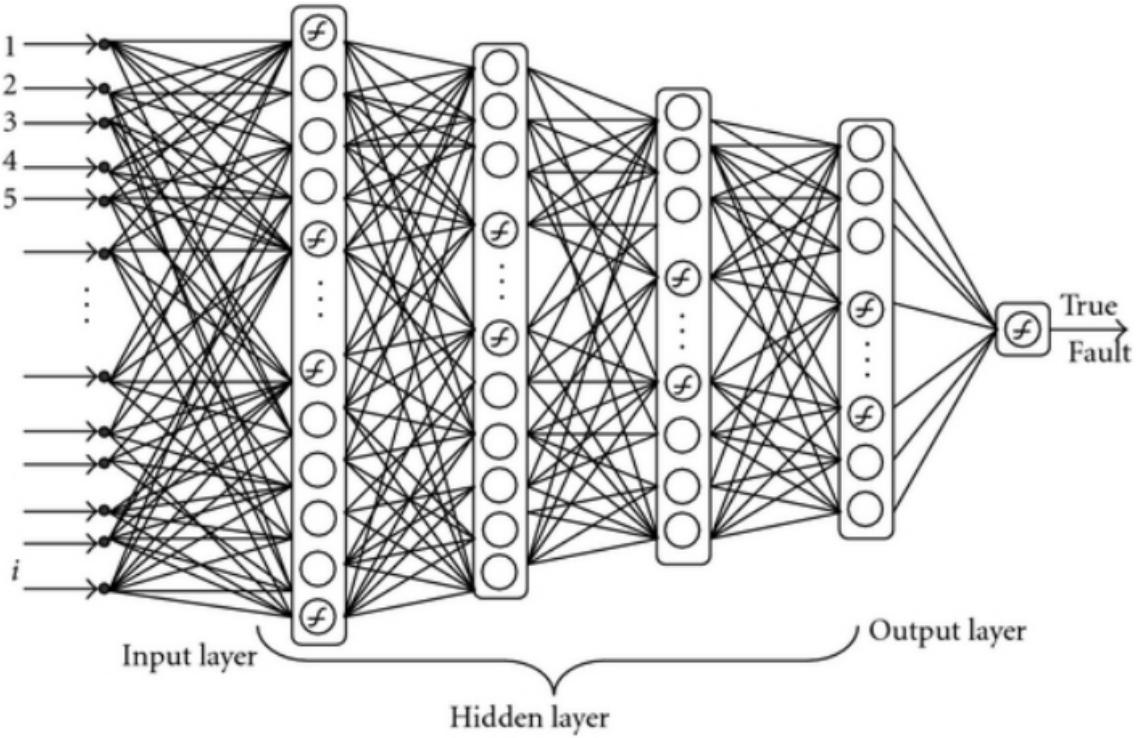
- o Which implies that:

$$t \leq \left(\frac{\alpha}{\beta}\right)^2 \quad (7)$$

- o Hence, if the perceptron made a  $t$ -th mistake,  $t$  must be less than or equal to  $(\frac{\alpha}{\beta})^2$ . That is, the perceptron always converges as long as the data set  $D$  is linearly separable (i.e. there exists a unit-length vector  $u$  s.t.  $y^{(i)}(u^T x^{(i)}) \geq \beta$  for  $\forall x^{(i)} \in D$ )
- o Limitations of Single Layer Perceptrons

- The representational power of the single-layer perceptron was not well understood at first in the AI community
  - In 1969, Minsky and Papert at MIT popularized a set of arguments showing that the single-layer perceptron could not learn certain classes of functions (including XOR)
  - They also showed that more complex functions could be represented using a multi-layer perceptron or MLP, but no one knew how to learn them from examples
  - This led to a shift away from mathematical/statistical models in AI toward logical/symbolic models

- Multi-Layer Perceptron



- The solution to MLP learning turned out to be:
  - Make the hidden layer non-linearities smooth (sigmoid/logistic) functions:

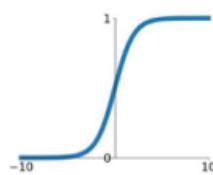
$$h_j^{(1)} = \frac{1}{1 + \exp(-(\sum_k w_{jk}^{(1)} x_k + b_j^{(1)}))} \quad (8)$$

$$h_j^{(t)} = \frac{1}{1 + \exp(-(\sum_k w_{jk}^{(t)} h_j^{(t-1)} + b_j^{(t)}))}$$

- Make the output layer non-linearity a smooth function
- Use standard numerical optimization methods (gradient descent) to learn the parameters. The algorithm is known as back propagation and was popularized by Rumelhart, Hinton and Willianms in the 1980s
- Activation Functions

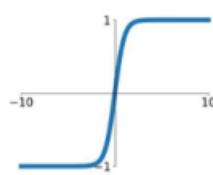
**Sigmoid**

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



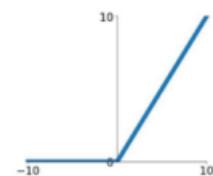
**tanh**

$$\tanh(x)$$



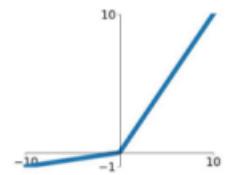
**ReLU**

$$\max(0, x)$$



**Leaky ReLU**

$$\max(0.1x, x)$$

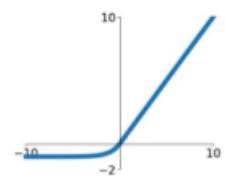


**Maxout**

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

**ELU**

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



- o Sigmoid or Logistic Activation Function

- Sigmoid function translates the input ranged in  $[-\infty; +\infty]$  to the range in  $(0, 1)$
- A more generalized sigmoid function that is used for **multiclass classification** called softmax function
- Problems with sigmoid function
  - The  $\exp(\cdot)$  function is computationally expensive
  - It has the problem of vanishing gradients
- Tanh Activation Function
  - The tanh function is defined as

$$\tanh(x) = \frac{e^{2x} - 1}{e^{2x} + 1} \quad (9)$$

- It is bound to range  $(-1, 1)$
- The gradient is stronger (i.e., steeper (陡)) for tanh than sigmoid
- Like sigmoid, tanh also has a vanishing gradient problem
- In practice, optimization is easier for tanh hence in practice it is always preferred over sigmoid function
- ReLU Activation Function
  - The ReLU (Rectified (矫正) Linear Units) function is defined as

$$ReLU(x) = \begin{cases} x, & \text{if } x \geq 0 \\ 0, & \text{if } x < 0 \end{cases} \quad (10)$$

- ReLU is identity for positive values, and zero for negative values
- Benefits of ReLU
  - Cheap to compute and easy to optimize
  - It converges faster
  - No vanishing gradient problem
  - Can output a true zero value, leading to representational sparsity
- Problems with ReLU: If one neuron gets negative it is unlikely for it to recover (since the gradient is zero). This is called "**dying ReLU**" problem
- Variants of ReLU
  - LReLU
    - Leaky ReLU is defined as:

$$LReLU(x) = \begin{cases} x, & \text{if } x \geq 0 \\ ax, & \text{if } x < 0 \end{cases} \quad (11)$$

- o Leaky ReLU attempts to fix the "dying ReLU" problem. Instead of the function being zero when  $x < 0$ , a leaky ReLU gives a small negative slope
- o  $a$  is a parameter constrained to be positive. It can either be pre-determined or learned from data.
- o ELU (Exponential Linear Unit)
  - The ELU is defined as:  $ELU(x) = \begin{cases} x, & \text{if } x \geq 0 \\ a(e^x - 1), & \text{if } x < 0 \end{cases}$

- It follows the same rule for  $x \geq 0$  as ReLU, and increases exponentially for  $x < 0$ .
- ELU tries to make the mean activations closer to zero which speeds up training (by adjusting  $a$ )
- Empirically, ELU leads to higher classification results
- Maxout Activation
  - The maxout activation function is defined as
$$\text{Maxout}(\mathbf{x}; \mathbf{w}_1, \mathbf{w}_2) = \max\{\mathbf{w}_1^T \mathbf{x}, \mathbf{w}_2^T \mathbf{x}\} \quad (12)$$
  - It is piecewise linear that returns the maximum of the inputs
  - The maxout activation is a generalization of ReLU and leaky ReLU. It is a learnable activation function
  - The maxout neuron, therefore, enjoys all the benefits of ReLU (linear regime of operation, no saturation) and does not have its drawbacks (dying ReLU)
  - However, it doubles the total number of parameters for each neuron and hence, a higher total number of parameters need to be trained
- Universal Approximation Theorem
  - It can be shown that a sigmoid network with one hidden layer (of infinite nodes) is a universal function approximator
    - In terms of classification, this means neural networks with one hidden layer (of unbounded size) can represent any decision boundary and thus have infinite capacity
    - It was also shown that deep networks can be more efficient at representing certain types of functions than shallow networks
    - It is not the specific choice of the activation function, but rather the multi-layer feedforward architecture itself which gives neural networks the potential of being universal approximators
    - It does not touch upon the algorithmic learnability of those parameters
  - More Failures
    - Through the 1990s, it became clear that while these models could represent arbitrarily complex functions and that deep networks should work better than shallow networks, no one could effectively train deep networks
    - This led to a shift away from neural networks towards statistical models and SVMs through the late 1990s and 2000s
    - Only a few groups continued to work on neural networks during this time period (Schmidhuber, Hinton, LeCun, Bengio, Ng)
    - Highly recommended reading: "Deep Learning: Our Miraculous Year 1990-1991" by Jürgen Schmidhuber

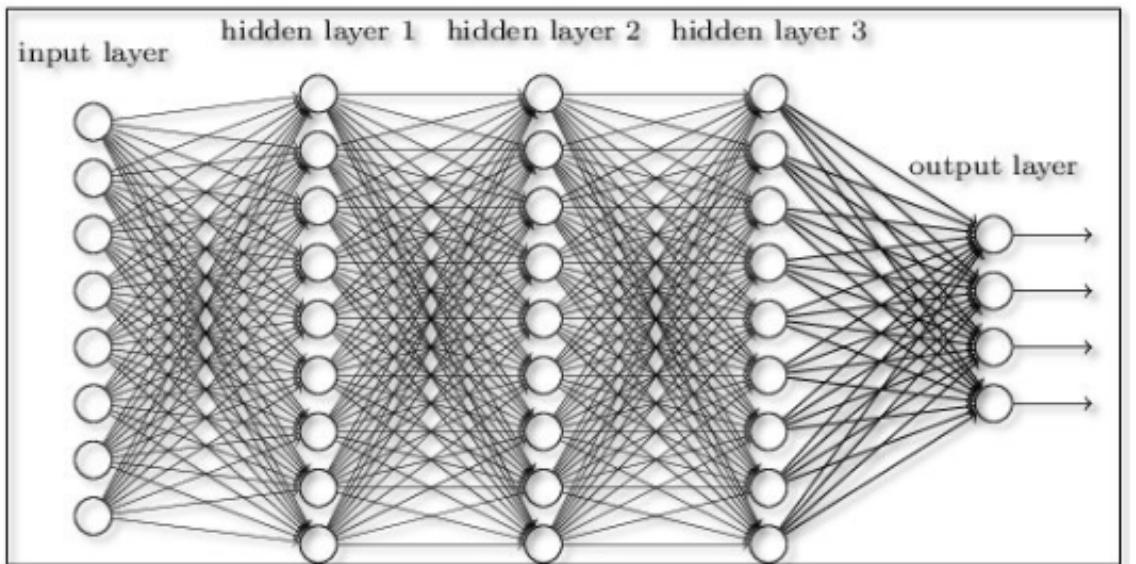
## Deep Learning

- Solution to the deep learning problem (as of right now) appears to be:
  - Have access to lots of labeled data (i.e., millions of examples)
  - Make the non-linearities non-smooth again (rectified linear units, ReLU):

$$h_j^{(1)} = \max(0, \sum_k w_{jk}^{(1)} x_k + b_j^{(1)}) \quad (13)$$

$$h_j^{(t)} = \max(0, \sum_k w_{jk}^{(t)} h_k^{(t-1)} + b_j^{(t)})$$

- Use improved optimization methods to learn the parameters
- Use new regularization schemes to constrain the large number of model parameters in the network
- Use GPUs for parallel computing with a massive 20-30 times speedup. Model training takes less than ten days for large vision problems instead of one year
- Main Idea
  - Learn a *feature hierarchy* all the way from raw data (i.e. image pixels for vision problems) to classifier
  - Each layer extracts features from the output of the previous layer
  - Train all layers jointly



## Backpropagation

- Backpropagation is a optimizing method to train deep neural network
- An example (non-vectorized)

## Backpropagation: a simple example

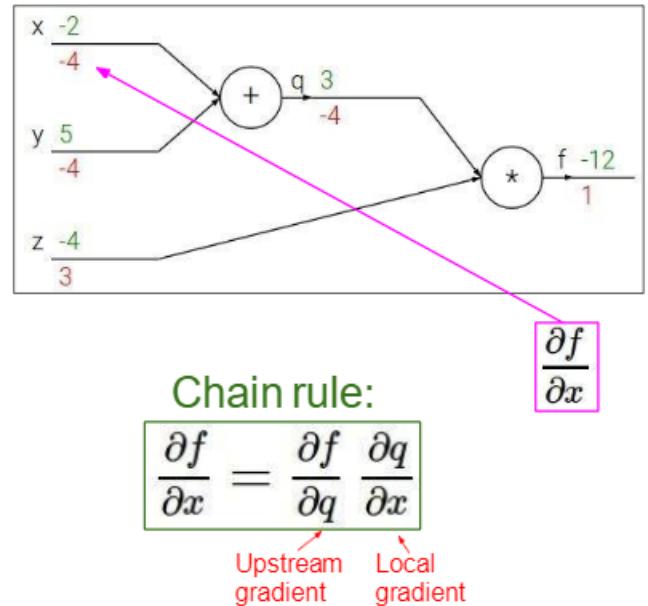
$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

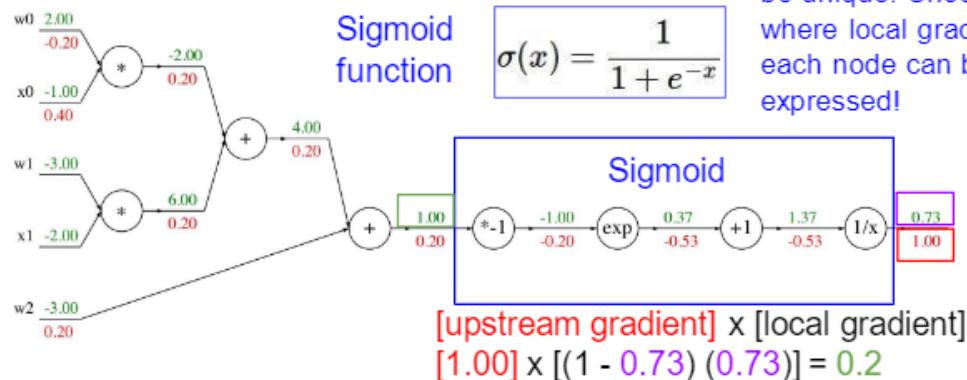
Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



- Another example:

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

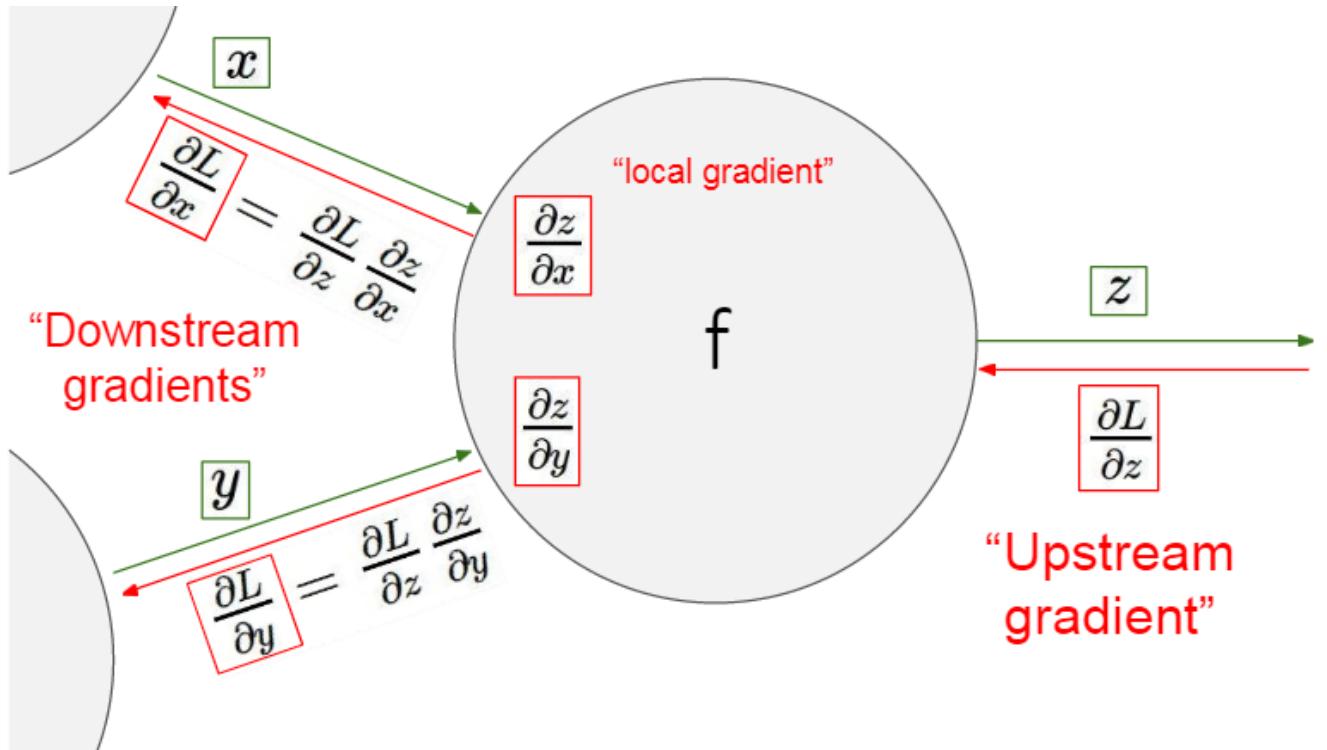


Computational graph representation may not be unique. Choose one where local gradients at each node can be easily expressed!

Sigmoid local gradient:

$$\frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = \left( \frac{1 + e^{-x} - 1}{1 + e^{-x}} \right) \left( \frac{1}{1 + e^{-x}} \right) = (1 - \sigma(x))\sigma(x)$$

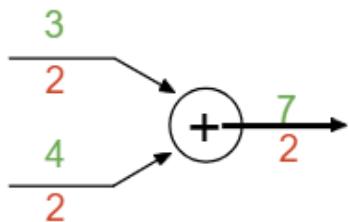
- General Case



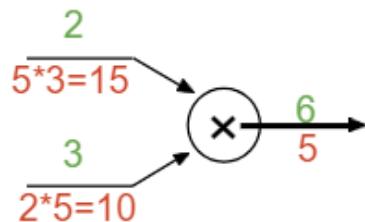
- Patterns

## Patterns in gradient flow

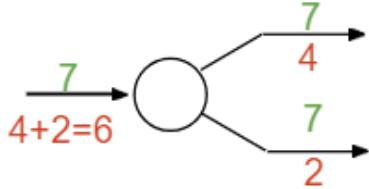
**add** gate: gradient distributor



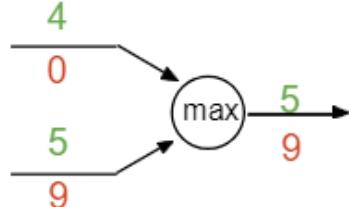
**mul** gate: “swap multiplier”



**copy** gate: gradient adder



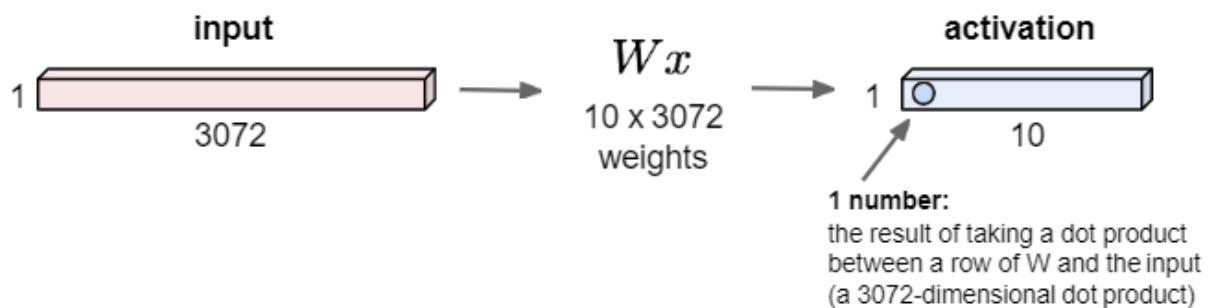
**max** gate: gradient router



## Convolutional Networks

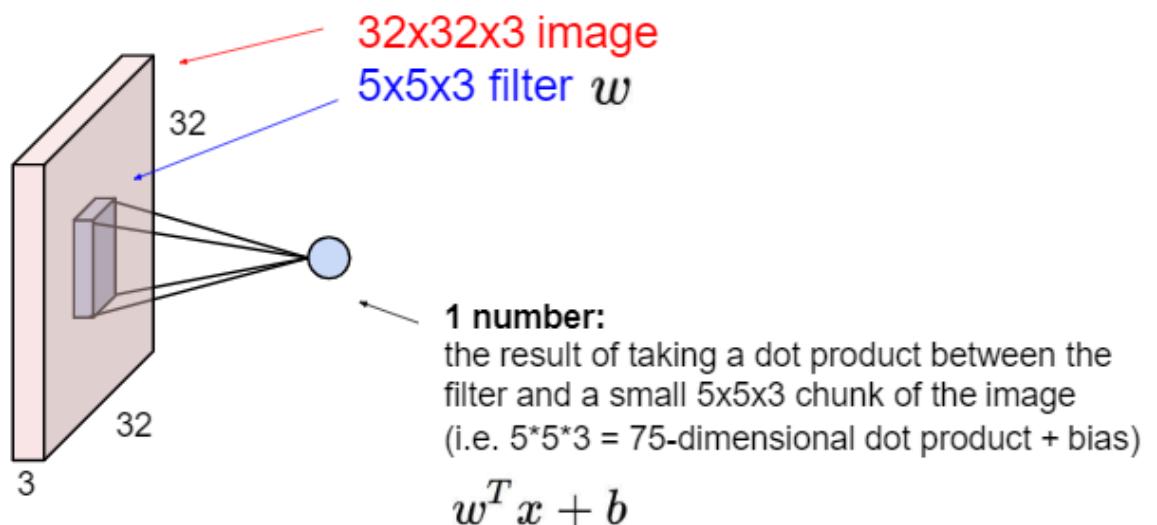
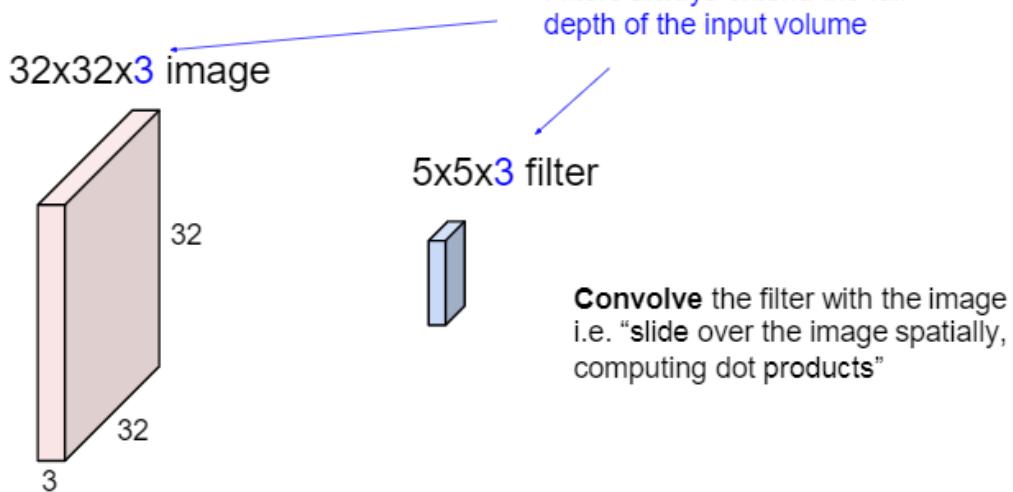
- Fully Connected Layer

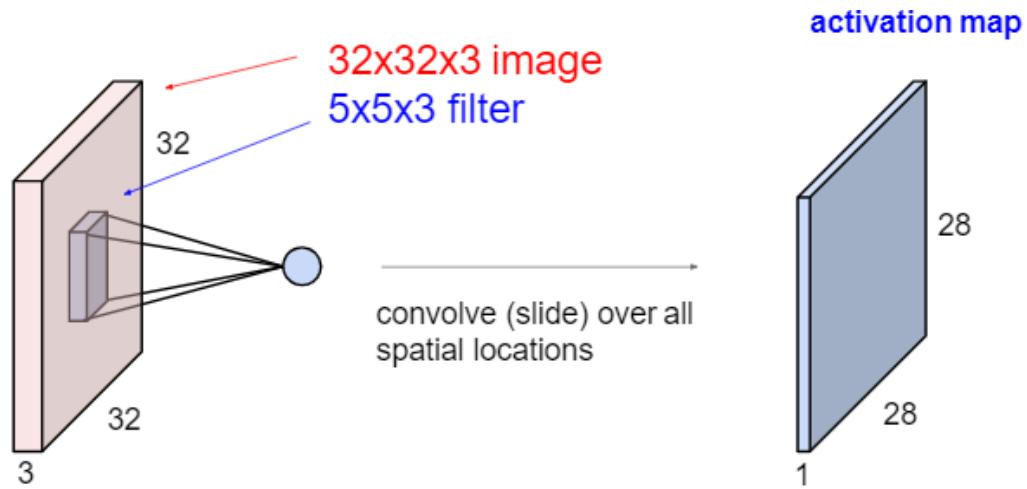
32x32x3 image  $\rightarrow$  stretch to 3072 x 1



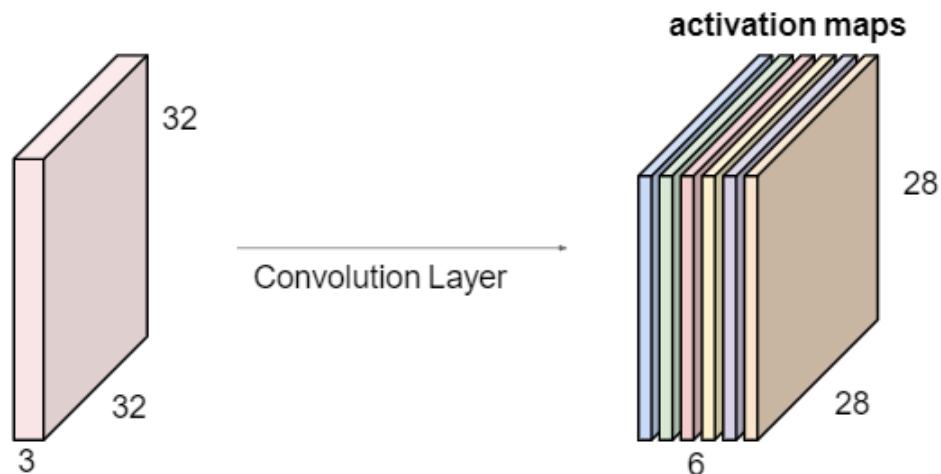
- Convolution Layer

## Convolution Layer



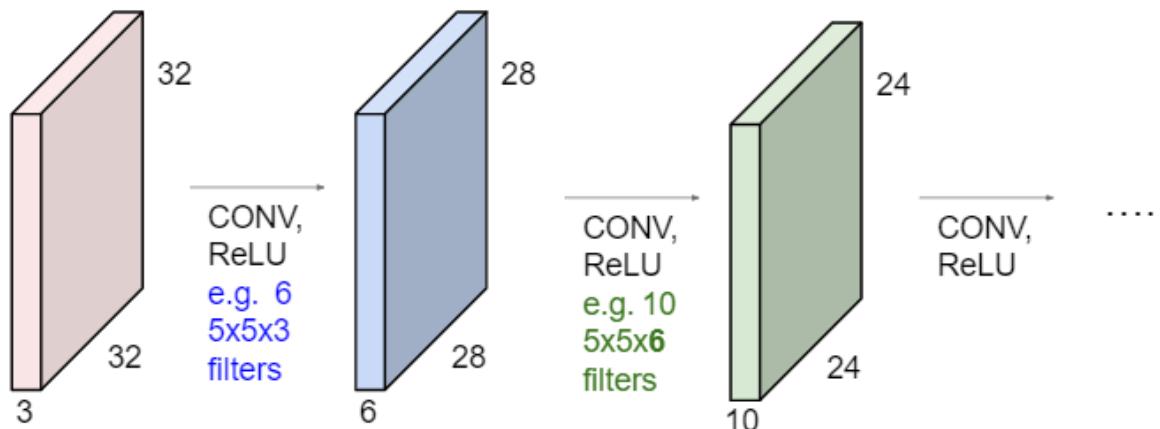


For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

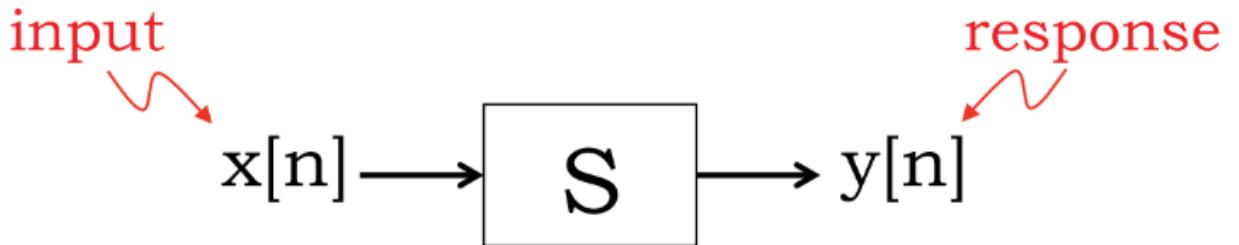


We stack these up to get a “new image” of size 28x28x6!

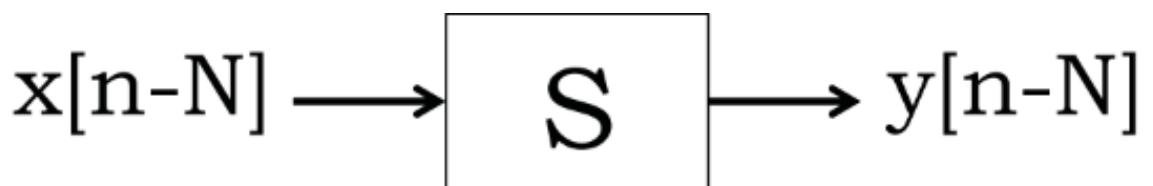
- ConvNet is a sequence of Convolution Layers, interspersed with activation function



- Origin of convolution idea
  - Convolution arises from the field of signal processing, which describes the output (in terms of the input) of an important class of operations (or systems) known as linear time-invariant(LTI)



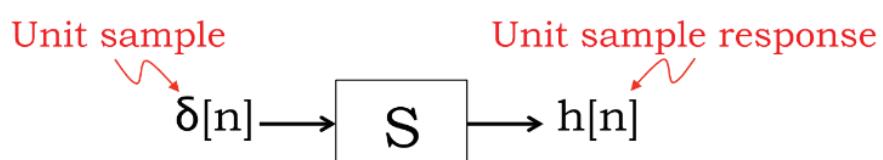
- A discrete-time signal such as  $x[n]$  or  $y[n]$  is described by an infinite sequence of values, *i.e.*, the time index  $n$  takes values in  $-\infty$  to  $\infty$ . The above picture is a snapshot at a particular time  $n$
- The sequence of output values  $y[.]$  is the response of system  $S$  to the input sequence  $x[.]$
- Time Invariant Systems
  - Let  $y[n]$  be the response of  $S$  to input  $x[n]$
  - If for all possible sequences  $x[n]$  and integers  $N$



- Then system  $S$  is said to be time invariant ( $TI$ ). A time shift in the input sequence to  $S$  results in an identical time shift of the output sequence
- Linear Systems
  - Let  $y_1[n]$  be the response of  $S$  to an arbitrary input  $x_1[n]$  and  $y_2[n]$  be the response to an arbitrary  $x_2[n]$
  - If, for arbitrary scalar coefficients  $a$  and  $b$ , we have

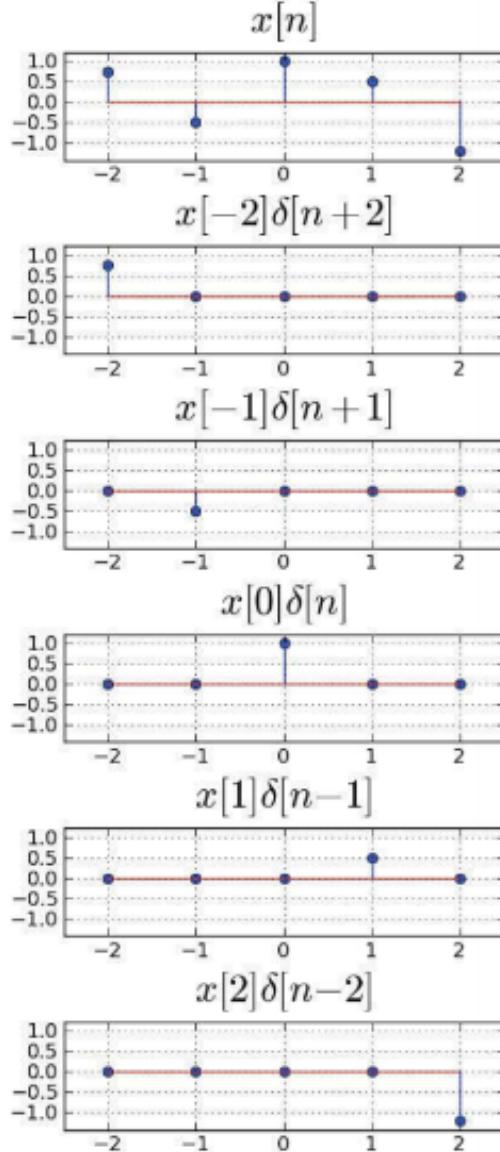


- Then system  $S$  is said to be linear. If the input is the weighted sum of several signals, the response is the superposition (*i.e.*, the same weighted sum) of the response to those signals
- One key consequence: If the input is identically 0 for a linear system, the output must also be identically 0
- Unit Sample Responses



$$\delta[n] = \begin{cases} 1, & \text{if } n = 0 \\ 0, & \text{otherwise} \end{cases} \quad \text{and} \quad \delta[n - N] = \begin{cases} 1, & \text{if } n = N \\ 0, & \text{otherwise.} \end{cases}$$

- The unit sample response of a system  $S$  is the response of the system to the unit sample input.  
We will always denote the unit sample response as  $h[n]$ .
- Unit Sample Decomposition
  - A discrete-time signal can be decomposed into a sum of time-shifted, scaled unit samples
  - Example:

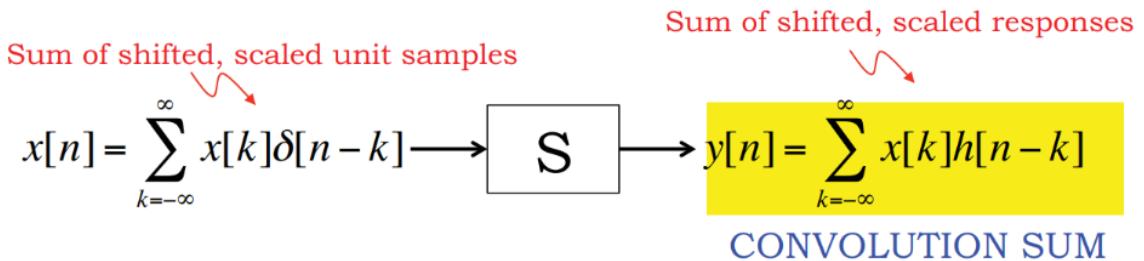


$$x[n] = x[-2]\delta[n+2] + x[-1]\delta[n+1] + x[0]\delta[n] + x[1]\delta[n-1] + x[2]\delta[n-2] \quad (14)$$

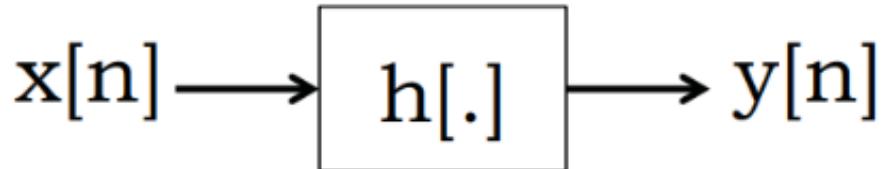
- In general, for any particular index, only one term of this sum is non-zero

$$x[n] = \sum_{k=-\infty}^{\infty} x[k]\delta[n-k] \quad (15)$$

- Modeling LTI Systems
  - If system  $S$  is both linear and time-invariant, then we can use the unit sample response to predict the response to any input waveform  $x[n]$



- Indeed, the unit sample response  $h[n]$  completely characterizes the LTI system  $S$ , so you often see



- Discrete Convolution: Discrete Convolution typically contains the following steps:

- List the index ' $k$ ' covering a sufficient range
- List the input  $x[k]$
- Obtain the reversed sequence  $h[-k]$ , and align the rightmost element of  $h[n - k]$  to the leftmost element of  $x[k]$
- Cross-multiply and sum (*i.e.*, dot product) the nonzero overlap terms to produce  $y[n]$
- Slide  $h[n - k]$  to the right by one position
- Repeat Steps 4 and 5; stop if all the output values are zero or if required

- Example

**Example:** Find the convolution of the two sequences  $x[n]$  and  $h[n]$  given by,

$$x[k] = [3 \ 1 \ 2] \quad h[k] = [3 \ 2 \ 1]$$

k:	-2	-1	0	1	2	3	4	5
x[k]:			3	1	2			
h[-k]:	1	2	3					
h[1-k]:		1	2	3				
h[2-k]:			1	2	3			
h[3-k]:				1	2	3		
h[4-k]:					1	2	3	
h[5-k]:						1	2	3

Hint: The value of  $k$  starts from  $(-\text{length of } h + 1)$  and continues till  $(\text{length of } h + \text{length of } x - 1)$

Here  $k$  starts from  $-3 + 1 = -2$  and continues till  $3 + 3 - 1 = 5$

$$y[0] = 3 \times 3 = 9$$

$$y[3] = 1 \times 1 + 2 \times 2 = 5$$

$$y[1] = 3 \times 2 + 3 \times 1 = 9$$

$$y[4] = 2 \times 1 = 2$$

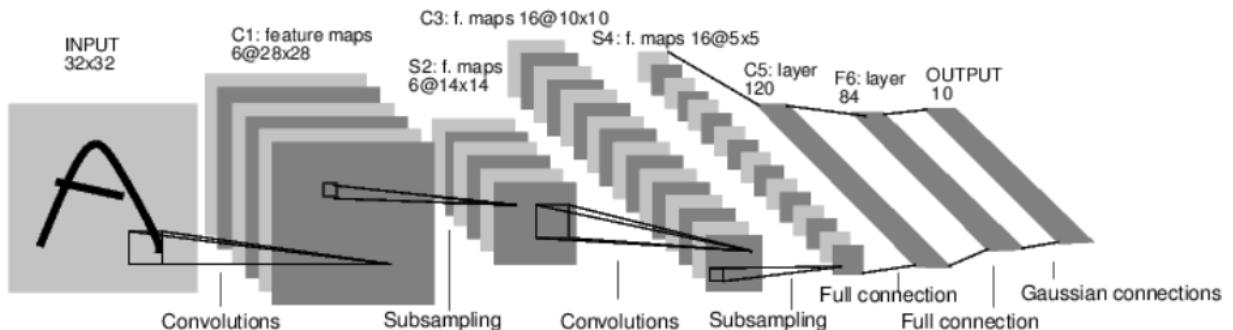
$$y[2] = 3 \times 1 + 1 \times 2 + 2 \times 3 = 11$$

$$y[5] = 0 \text{ (no overlap)}$$

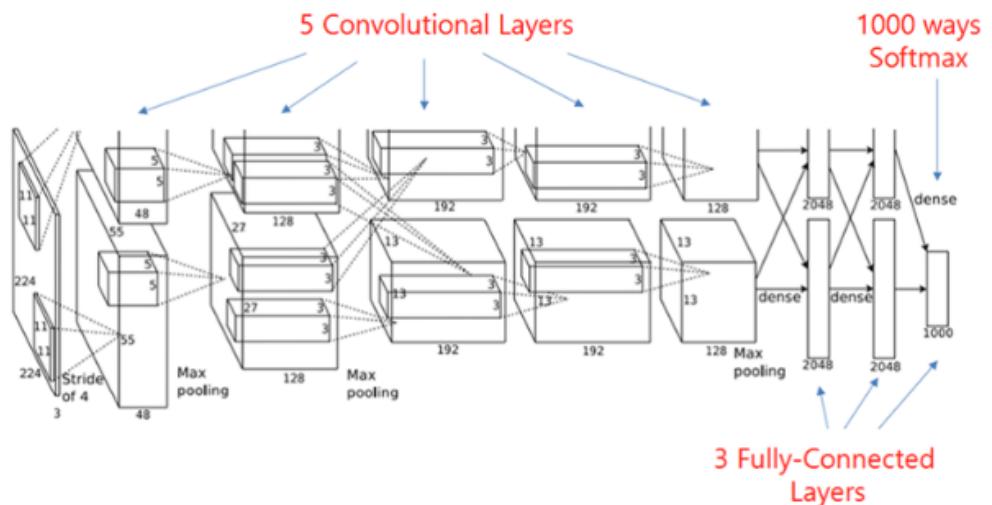
$$y[n] = \boxed{\{9 \ 9 \ 11 \ 5 \ 2 \ 0\}}$$

- Representative CNNs

- LeNet-5



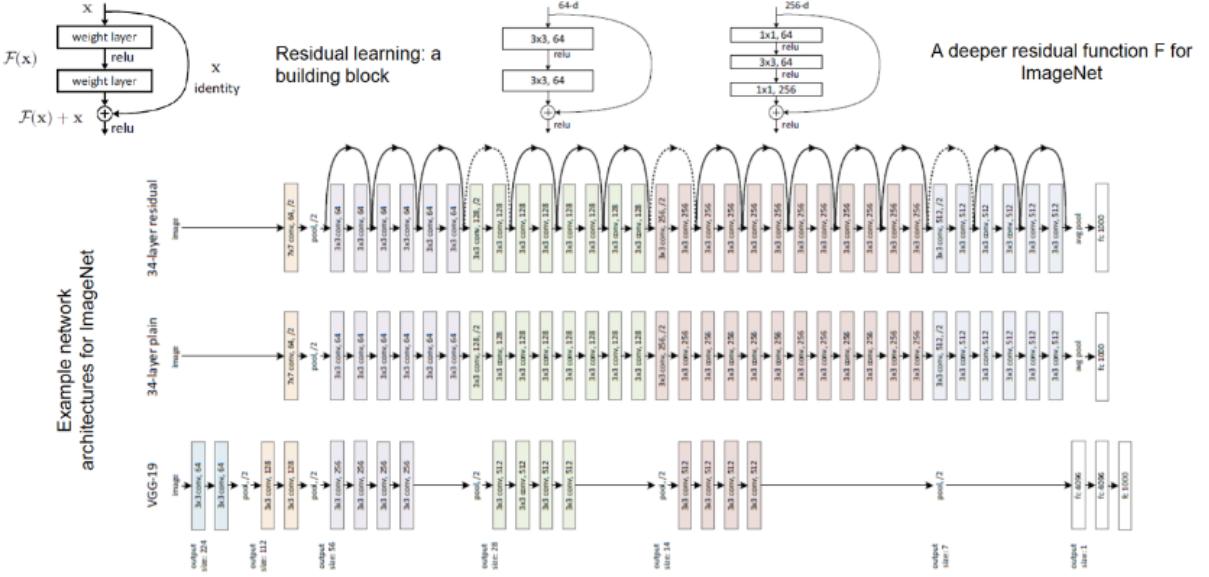
- Average pooling (sub-sampling)
  - Sigmoid or tanh nonlinearity
  - Fully connected layers at the end
  - Trained on MNIST digit dataset with 60K training examples
  - Ref: Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," Proc. IEEE, 1998.
- AlexNet, 2012



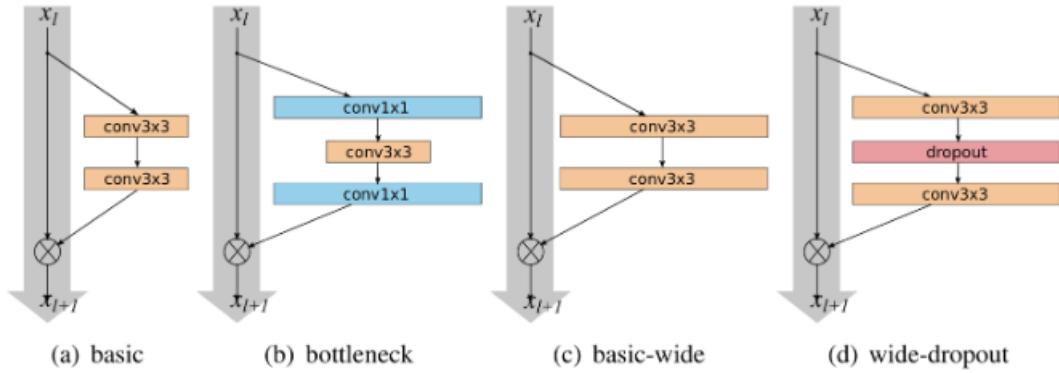
- The **first** winner deep model in computer vision
  - 5 convolutional layers + 3 fully-connected layers + softmax classifier
  - Key technical features: ReLU, dropout, data augmentation
  - Ref: A. Krizhevsky, I. Sutskever, and G. Hinton, "Imagenet classification with deep convolutional neural networks." inNIPS, 2012.
- VGG-Net, 2014

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input ( $224 \times 224$ RGB image)					
conv3-64	conv3-64 LRN	conv3-64 <b>conv3-64</b>	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 <b>conv3-128</b>	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256	conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 <b>conv1-256</b>	conv3-256 conv3-256 <b>conv3-256</b>
maxpool					
conv3-512	conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
conv3-512	conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

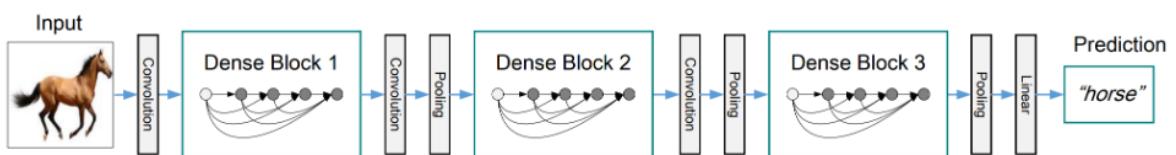
- Key technical features: Increased depth (up to 19) and smaller filter size (3)
- Configurations D and E are widely used for various tasks, called VGG-16 and VGG-19
- Ref: K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in ICLR, 2015.
- Deep Residual Network (ResNet), 2015

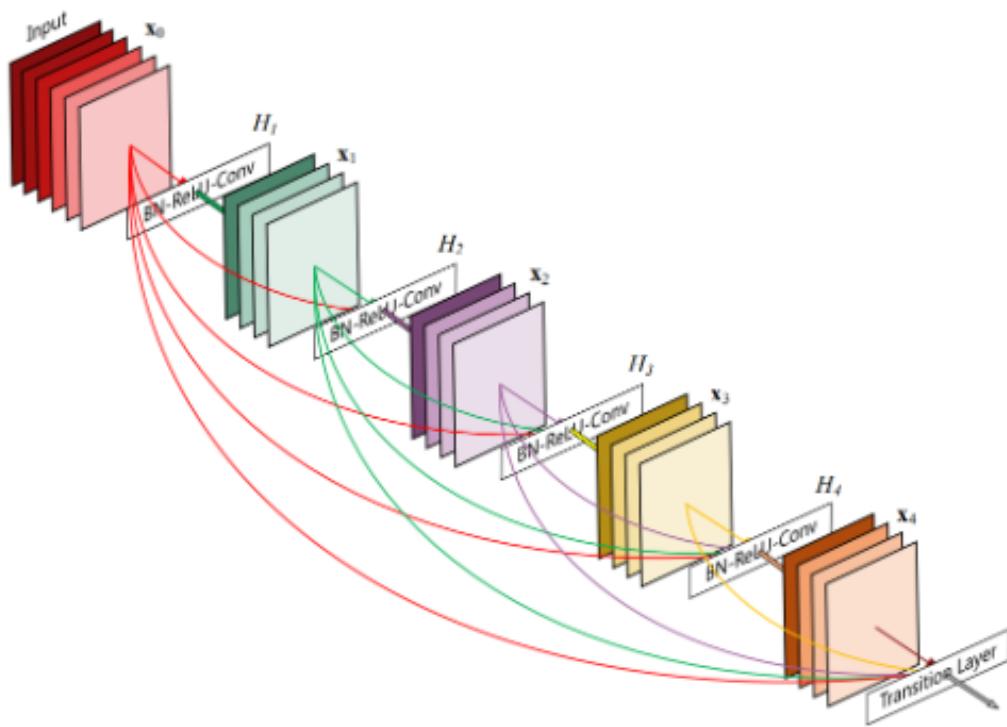


- Key technical features: Skip connections for residual mapping up to  $> 1000$  layers
- Ref: K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," inCVPR, 2016.
- Wide ResNet, 2016

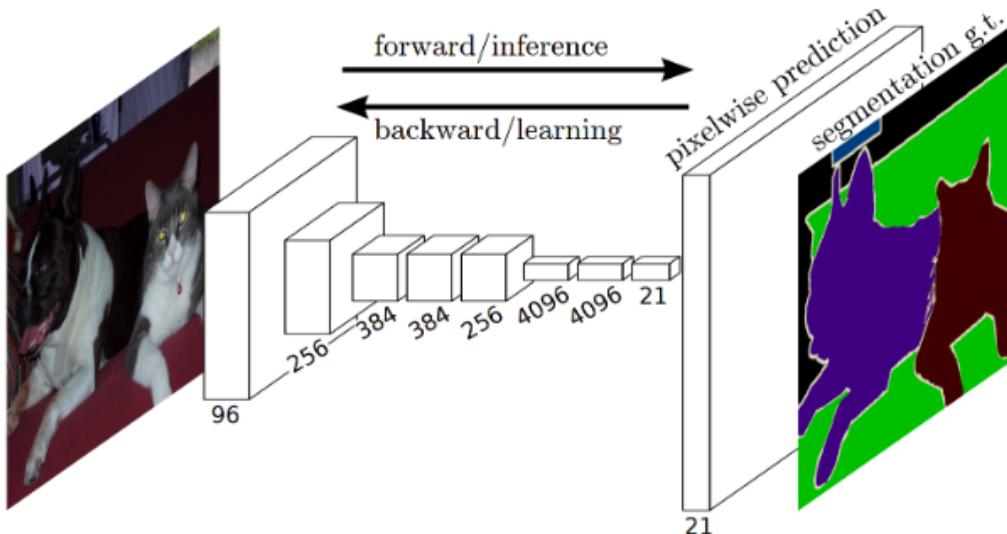


- Widening of ResNet blocks (if done properly) provides a more effective way of improving performance of residual networks compared to increasing their depth
- A wide 16-layer deep network has the same accuracy as a 1000-layer thin deep network and a comparable number of parameters, although being several times faster to train
- Ref: S. Zagoruyko and N. Komodakis, "Wide residual networks," inBMVC, 2016.
- Densely Connected Networks (DenseNet), 2017



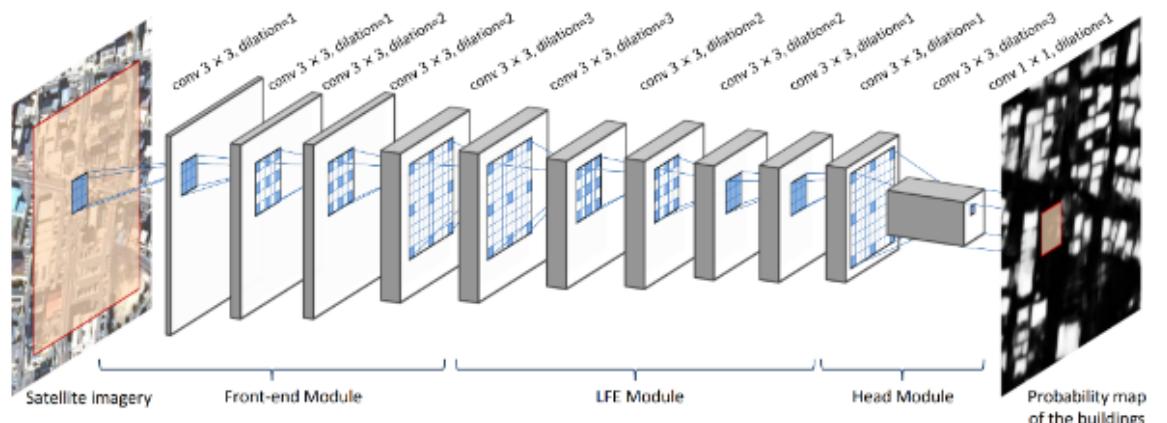
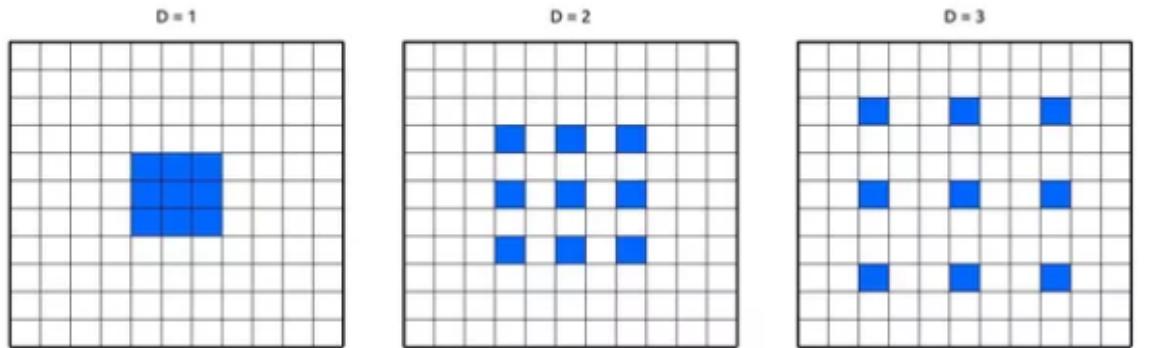


- Key technical features: Finer combination of multi-scale features
- Ref: G. Huang, Z. Liu, L. van der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in CVPR, 2017.
- More Art of Convolutions:
  - Fully Convolutional Networks (FCN), 2014



- Key technical features: No fully-connected layer => No fixed requirement on input size
- Widely adopted in pixel-to-pixel prediction tasks, e.g., image segmentation
- Ref: J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in CVPR, 2015.
- U-Net, 2015
  - The architecture consists of a contracting path to capture context
  - and a symmetric expanding path to enable precise localization

- Very popular backbone for dense prediction (image segmentation, restoration, etc.)
- Ref: O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional networks for biomedical image segmentation," in MICAI, 2015.
- Dilated Convolution, 2015



- Aggregates multi-scale contextual information without losing resolution or analyzing rescaled images
- Has layer receptive fields with the same computation and memory cost
- Improves the performance on dense prediction tasks, e.g., semantic segmentation
- Ref: F. Yu, and K. Vladlen, "Multi-scale context aggregation by dilated convolutions," in ICLR, 2016.

## Lecture 13: Deep Learning -- Optimization

### Optimization and Implementation

- Gradient Descent (GD)

---

## Algorithm 1 Batch Gradient Descent at Iteration $t$

---

**Require:** Learning Rate  $\alpha^{(t)}$

**Require:** Initial Parameter  $\theta^{(0)}$

- 1: **while** stopping criteria not met **do**
  - 2:   Compute gradient estimate over  $m$  examples:
  - 3:    $\hat{\mathbf{g}} = \frac{1}{m} \nabla_{\theta^{(t)}} \sum_{i=1}^m \ell(f(\mathbf{x}^{(i)}; \theta^{(t)}), y^{(i)})$
  - 4:   Apply update:  $\theta^{(t+1)} = \theta^{(t)} - \alpha^{(t)} \hat{\mathbf{g}}$
  - 5:    $t \leftarrow t + 1$
  - 6: **end while**
- 

- Pros: Gradient estimates are stable
- Cons: Need to compute gradients over the entire training set for one update
- Stochastic Gradient Descent (SGD)

---

## Algorithm 2 Stochastic Gradient Descent at Iteration $t$

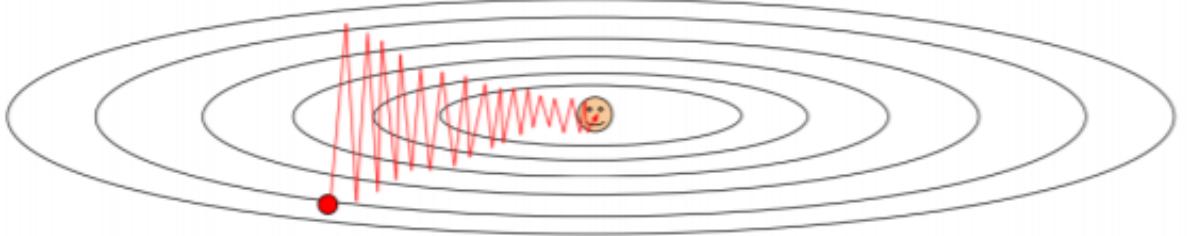
---

**Require:** Learning Rate  $\alpha^{(t)}$

**Require:** Initial Parameter  $\theta^{(0)}$

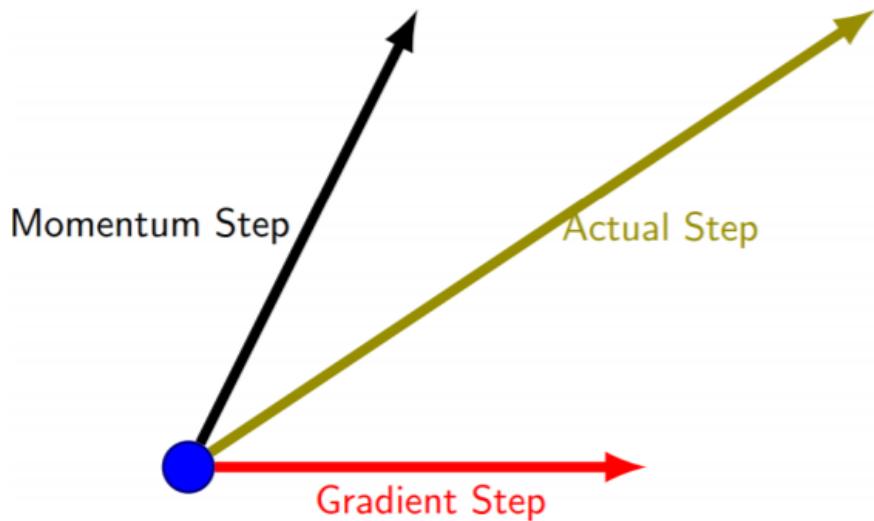
- 1: **while** stopping criteria not met **do**
  - 2:   Sample example  $(\mathbf{x}^{(i)}, y^{(i)})$  from the training set
  - 3:   Compute gradient estimate:
  - 4:    $\hat{\mathbf{g}} = \nabla_{\theta^{(t)}} \ell(f(\mathbf{x}^{(i)}; \theta^{(t)}), y^{(i)})$
  - 5:   Apply update:  $\theta^{(t+1)} = \theta^{(t)} - \alpha^{(t)} \hat{\mathbf{g}}$
  - 6:    $t \leftarrow t + 1$
  - 7: **end while**
- 

- Pros: Computation time per update does not depend on  $m$ , allowing convergence on extremely large datasets
- Cons: Gradient estimates can be noisy. Use mini batches
- Problems with GD and SGD
  - What if loss changes quickly in one direction and slowly in another (*i. e.*, error surface has high curvature)



- What does gradient descent do: Very slow progress along shallow dimension, jittering along steep direction
- Momentum
  - How do we try and solve this problem
  - Introduce a new variable  $\mathbf{v}$ , the velocity
  - We think of  $\mathbf{v}$  as the direction and speed by which the parameters move as the learning dynamics progress
  - The velocity is an exponentially decaying moving average of the negative gradients

$$\begin{aligned}\mathbf{v}^{t+1} &= \rho \mathbf{v}^t - \alpha^{(t)} \nabla_{\theta^{(t)}} l(f(\mathbf{x}^{(i)}; \theta^{(t)}), y^{(i)}) \\ \alpha &\in [0, 1) \\ \theta^{(t+1)} &= \theta^{(t)} + \mathbf{v}^{(t+1)}\end{aligned}\tag{16}$$



- The velocity accumulates the previous gradients
- The role of  $\rho$ 
  - If  $\rho$  is larger than  $\alpha^{(t)}$ , the current update is more affected by the previous gradients
  - Usually values for  $\rho$  are set high  $\approx 0.9$
- Performance

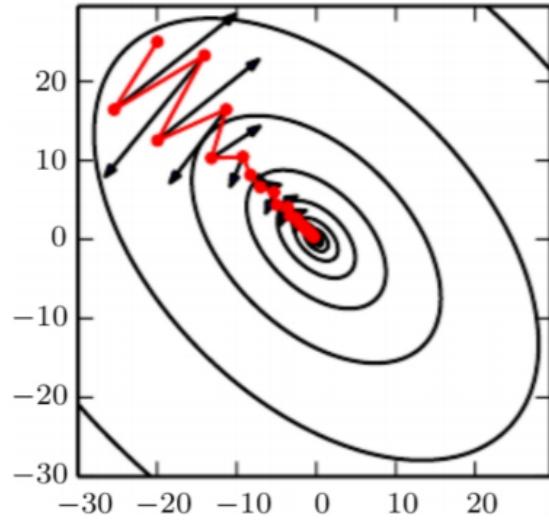


Illustration of how momentum traverses such an error surface better compared to SGD

- Pseudo code

---

### **Algorithm 3** SGD with Momentum at Iteration $t$

---

**Require:** Learning Rate  $\alpha^{(t)}$

**Require:** Momentum Parameter  $\rho$

**Require:** Initial Parameter  $\theta^{(0)}$

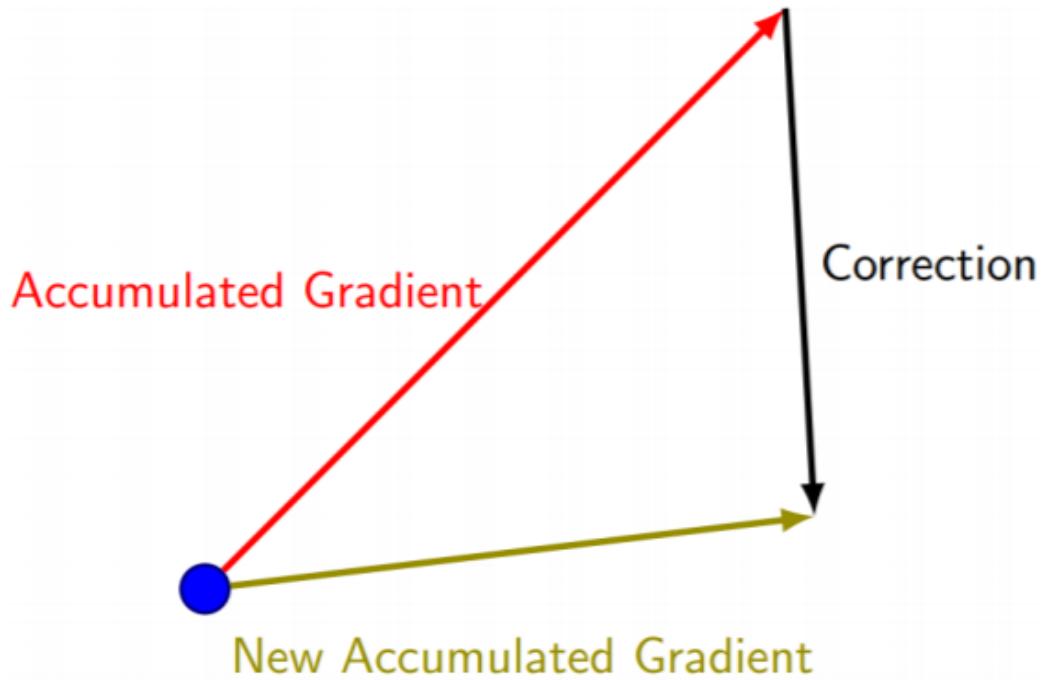
**Require:** Initial Velocity  $v^{(0)}$

- 1: **while** stopping criteria not met **do**
  - 2:    Sample example  $(\mathbf{x}^{(i)}, y^{(i)})$  from the training set
  - 3:    Compute gradient estimate:  $\hat{\mathbf{g}} = \nabla_{\theta^{(t)}} \ell(f(\mathbf{x}^{(i)}; \theta^{(t)}), y^{(i)})$
  - 4:    Compute the velocity update:  $\mathbf{v}^{(t+1)} = \rho \mathbf{v}^{(t)} - \alpha^{(t)} \hat{\mathbf{g}}$
  - 5:    Apply update:  $\theta^{(t+1)} = \theta^{(t)} + \mathbf{v}^{(t+1)}$
  - 6:     $t \leftarrow t + 1$
  - 7: **end while**
- 

- Nesterov Momentum

- Another approach:

- First take a step in the direction of the accumulated gradient
- Then calculate the gradient and make a correction



- Equation

$$\begin{aligned}\hat{\theta} &= \theta^{(t)} + \rho v^{(t)} \\ v^{(t+1)} &= \rho v^t - \alpha^{(t)} \nabla_{\theta^{(t)}} l(f(\mathbf{x}^{(i)}; \hat{\theta}), y^{(i)}) \\ \theta^{(t+1)} &= \theta^{(t)} + v^{(t+1)}\end{aligned}\tag{17}$$

- Pseudo code

---

#### **Algorithm 4** SGD with Nesterov Momentum at Iteration $t$

---

**Require:** Learning Rate  $\alpha^{(t)}$

**Require:** Momentum Parameter  $\rho$

**Require:** Initial Parameter  $\theta^{(0)}$

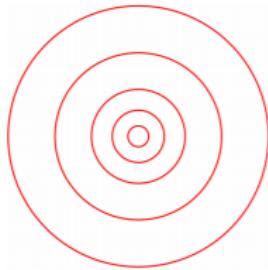
**Require:** Initial Velocity  $v^{(0)}$

- 1: **while** stopping criteria not met **do**
  - 2:    Sample example  $(\mathbf{x}^{(i)}, y^{(i)})$  from the training set
  - 3:    Update parameters:  $\tilde{\theta} = \theta^{(t)} + \rho v^{(t)}$
  - 4:    Compute gradient estimate:  $\hat{\mathbf{g}} = \nabla_{\tilde{\theta}} l(f(\mathbf{x}^{(i)}; \tilde{\theta}), y^{(i)})$
  - 5:    Compute the velocity update:  $v^{(t+1)} = \rho v^{(t)} - \alpha^{(t)} \hat{\mathbf{g}}$
  - 6:    Apply update:  $\theta^{(t+1)} = \theta^{(t)} + v^{(t+1)}$
  - 7:     $t \leftarrow t + 1$
  - 8: **end while**
- 

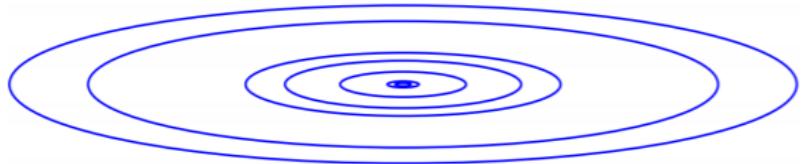
- Adaptive Learning Rate Methods

- Till now we assign the same learning rate to all parameters  $\theta_j$ 's

- If  $\theta_j$ 's vary in importance and frequency, why is this a good idea?
- It's probably not



Nice (all features are  
equally important)



Harder!

- AdaGrad

---

## Algorithm 5 AdaGrad

---

**Require:** Global Learning Rate  $\alpha$

**Require:** Initial Parameter  $\theta^{(0)}$

- 1: Initialize  $\mathbf{r}^{(0)} = 0$
  - 2: **while** stopping criteria not met **do**
  - 3:     Sample example  $(\mathbf{x}^{(i)}, y^{(i)})$  from the training set
  - 4:     Compute gradient estimate:  $\hat{\mathbf{g}} = \nabla_{\theta^{(t)}} \ell(f(\mathbf{x}^{(i)}; \theta^{(t)}), y^{(i)})$
  - 5:     Accumulate:  $\mathbf{r}^{(t+1)} = \mathbf{r}^{(t)} + \hat{\mathbf{g}} \odot \hat{\mathbf{g}}$
  - 6:     Compute update:  $\Delta\theta = -\frac{\alpha}{\delta + \sqrt{\mathbf{r}^{(t+1)}}} \odot \hat{\mathbf{g}}$
  - 7:     Apply update:  $\theta^{(t+1)} = \theta^{(t)} + \Delta\theta$
  - 8:      $t \leftarrow t + 1$
  - 9: **end while**
- 

- $\delta$  is used to make sure not dividing zero
- Idea: Scale the gradient of a model parameter by square root of sum of squares of all its historical values
- Progress along “steep” directions (with large partial derivatives) is damped
- Progress along “flat” directions (with small partial derivatives) is accelerated
- What happens to the step size over long time? Decays to zero
- RMSProp: “Leaky AdaGrad”
  - AdaGrad is good when the objective is convex
  - AdaGrad might shrink the learning rate too aggressively, we want to keep the history in mind
  - We can adapt it to perform better in non-convex settings by accumulating an exponentially decaying average of gradient

- This is an idea that we use again and again in deep learning
- Has about 2,500 citations on Google Scholar (As of Aug. 2019), but was proposed in a slide Geoffrey Hinton’s Coursera course

## Algorithm 6 RMSProp

**Require:** Global Learning Rate  $\alpha$

**Require:** Decay Parameter  $\rho$

**Require:** Initial Parameter  $\theta^{(0)}$

1: Initialize  $\mathbf{r}^{(0)} = 0$

2: **while** stopping criteria not met **do**

3:    Sample example  $(\mathbf{x}^{(i)}, y^{(i)})$  from the training set

4:    Compute gradient estimate:  $\hat{\mathbf{g}} = \nabla_{\theta^{(t)}} \ell(f(\mathbf{x}^{(i)}; \theta^{(t)}), y^{(i)})$

5:    Accumulate:  $\mathbf{r}^{(t+1)} = \rho \mathbf{r}^{(t)} + (1 - \rho) \hat{\mathbf{g}} \odot \hat{\mathbf{g}}$

6:    Compute update:  $\Delta\theta = -\frac{\alpha}{\delta + \sqrt{\mathbf{r}^{(t+1)}}} \odot \hat{\mathbf{g}}$

7:    Apply update:  $\theta^{(t+1)} = \theta^{(t)} + \Delta\theta$

8:     $t \leftarrow t + 1$

9: **end while**

- Adam: ADaptive Moments
  - Adam is currently THE default optimization algorithm for training deep neural networks
  - Introduced by Kingma and Ba in 2015 and quickly accumulated more than 25,000 Google Scholar Citations as of Aug. 2019
  - Adam is like RMSProp with Momentum but with bias correction terms for the first and second moments

---

## Algorithm 7 Adam

---

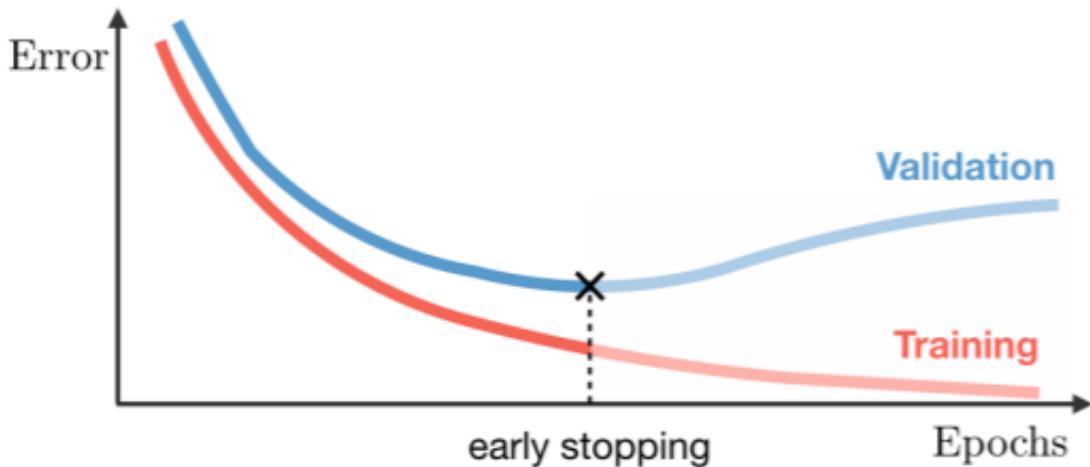
**Require:** Global Learning Rate:  $\alpha$  (set to 0.0001), Decay Rates:  $\rho_1$  (set to 0.9),  $\rho_2$  (set to 0.9), and Initial Parameter:  $\theta^{(0)}$

- 1: Initialize moment variables  $s^{(0)} = 0$  and  $r^{(0)} = 0$
- 2: **while** stopping criteria not met **do**
- 3:   Sample example  $(\mathbf{x}^{(i)}, y^{(i)})$  from the training set
- 4:   Compute gradient estimate:  $\hat{\mathbf{g}} = \nabla_{\theta^{(t)}} \ell(f(\mathbf{x}^{(i)}; \theta^{(t)}), y^{(i)})$
- 5:   Update:  $s^{(t+1)} = \rho_1 s^{(t)} + (1 - \rho_1) \hat{\mathbf{g}}$
- 6:   Update:  $r^{(t+1)} = \rho_2 r^{(t)} + (1 - \rho_2) \hat{\mathbf{g}} \odot \hat{\mathbf{g}}$
- 7:   Correct biases:  $\hat{s} = \frac{s^{(t+1)}}{1 - \rho_1^{t+1}}$  and  $\hat{r} = \frac{r^{(t+1)}}{1 - \rho_2^{t+1}}$
- 8:   Compute and apply update:  $\Delta \theta = -\alpha \frac{\hat{s}}{\hat{s} + \sqrt{\hat{r}}}$ ,  $\theta^{(t+1)} = \theta^{(t)} + \Delta \theta$
- 9:    $t \leftarrow t + 1$
- 10: **end while**

---

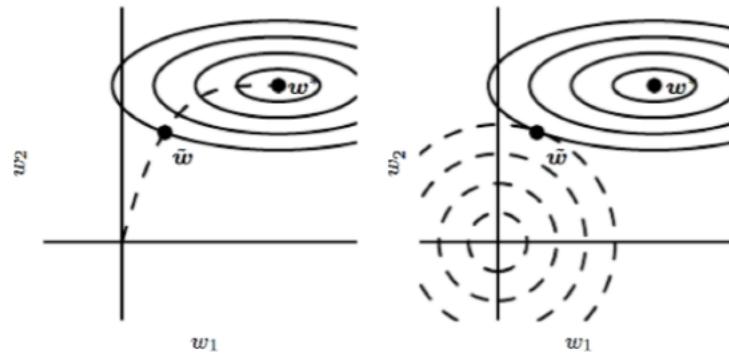
## Regularizations

- Data Augmentation
  - Adding noise to the input: a special kind of augmentation
  - Be careful about the transformation applied => label preserving
    - Example: classifying 'b' and 'd'; '6' and '9'
- Early Stopping



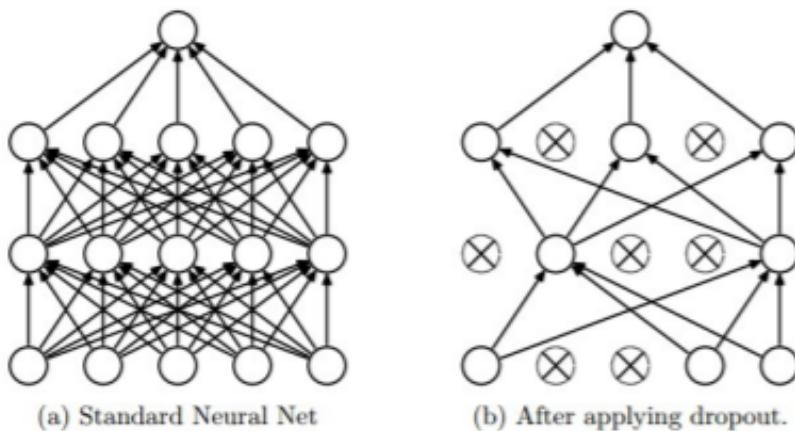
- Idea: Don't train the network to small training error
- Recall overfitting: With a larger hypothesis class, it is easier to find a hypothesis that fits training samples
- Prevent overfitting: Use validation error to decide when to stop

## Early Stopping as Regularization



Qualitatively, the number of degrees of freedom in the network starts out small and grows during the training process, corresponding to a steady increase in the effective complexity of the model. Early stopping then represents a way of limiting the effective network complexity

- In practice, when training, also output validation error
- Every time validation error improves, store a copy of the network weights
- When validation error plateaus for some time, stop
- Return the copy of the weights stored
- The number of training steps (i.e., iterations) is also a hyperparameter
- Dropout



- In each update step, randomly sample a different binary mask to all the input and hidden units
- Multiple the mask with the units and do the update as usual
- Typical dropout probability: 0.8 for input and 0.5 for hidden units
- Very useful for fully connected layers, less for convolutional layers
- How can this possibly be a good idea?

- Forces the network to have a redundant representation; prevents co-adaptation of features
  - Trains a large ensemble of models (that share parameters); Each binary mask corresponds to one model
  - How to vote to make a final prediction
    - Mask sampling: Randomly sample some (typically, 10 ~20) masks
    - For each mask, apply it to the trained model to get a prediction
    - Average the prediction
    - Ref: N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A simple way to prevent neural networks from overfitting," JMLR, 2014.
- Batch Normalization
  - Idea: Adjusts activations to lie within a desired operating range, while maintaining their relative values
  - Given a mini-batch  $B = \{(\mathbf{x}^{(i)}, y^{(i)}), i = 1, \dots, m\}$
  - Compute per-channel (*i.e.* feature) mean:  $\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$
  - Compute per-channel variance:  $\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$
  - Normalize  $\mathbf{x}^{(i)}$  across channel:  $\hat{x}_j^{(i)} = \frac{x_j^{(i)} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$
  - Issue: What if zero-mean, unit variance is too hard of a constraint
  - Incorporate learnable scale and shift parameters  $\gamma$  and  $\beta$ 

$$z_j^{(i)} = \gamma_j \hat{x}_j^{(i)} + \beta_j \quad (18)$$
  - Test time
    - Mean and variance estimates depend on mini-batch; can't do this at test-time
    - Solution: Keep exponentially decaying moving average of values seen during training and use them for testing
    - During testing, batch normalization becomes a linear operator; can be fused with the previous convolutional layer
    - Learning  $\gamma_j = \sigma_j$  and  $\beta_j = \mu_j$  will recover the identity function
    - Usually inserted after convolutional layers and before nonlinearity
    - Behaves differently during training and testing: this is a very common source of bugs

## Practical Tricks

- Data pre-processing
  - Simplest: Zero-center the data, and then normalize them
    - Makes sense only when data attributes have approximately equal importance but different scales
  - PCA whitening: The data is first centered, and then projected into the eigenbasis, followed by dividing every dimension by the corresponding eigenvalue
  - Never forget to shuffle your data each epoch for SGD input
- Weight Initialization
  - All zero initialization: Terribly Wrong!

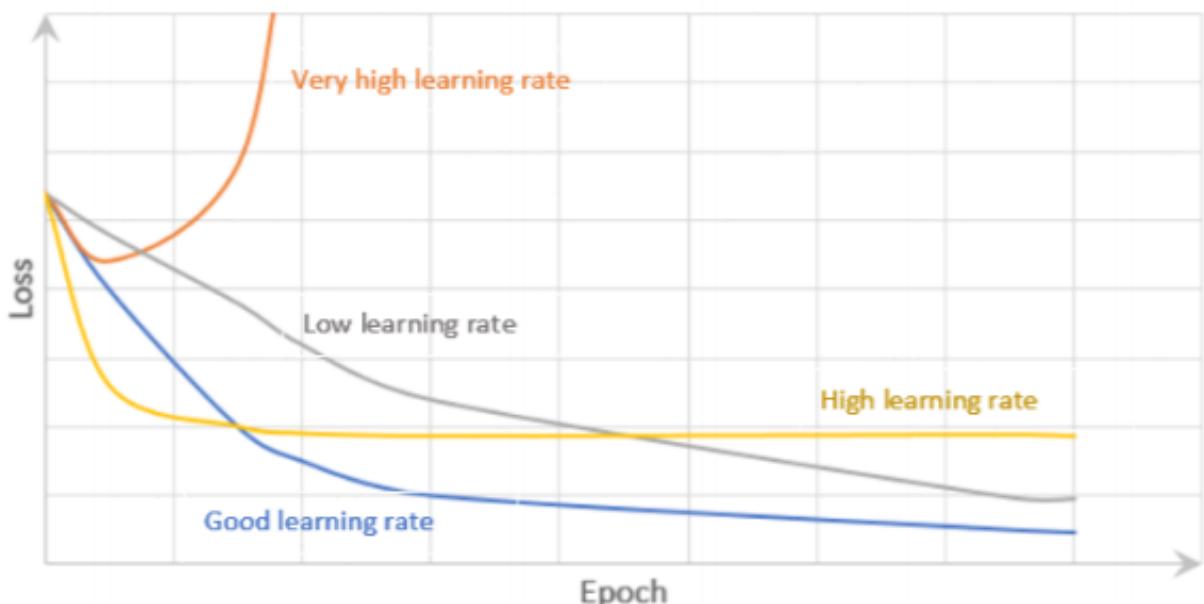
- If every neuron in the network computes the same output, then all of them will also compute the same gradients during backpropagation and undergo the exact same parameter updates
- Small random initialization
  - Works okay for small networks, but not for deeper networks
- Current recommendation for initializing CNNs with ReLU:

```
1 | w = np.random.randn(n) *sqrt(2.0/n)
```

'randn': Gaussian; 'n': the number of inputs for the current layer

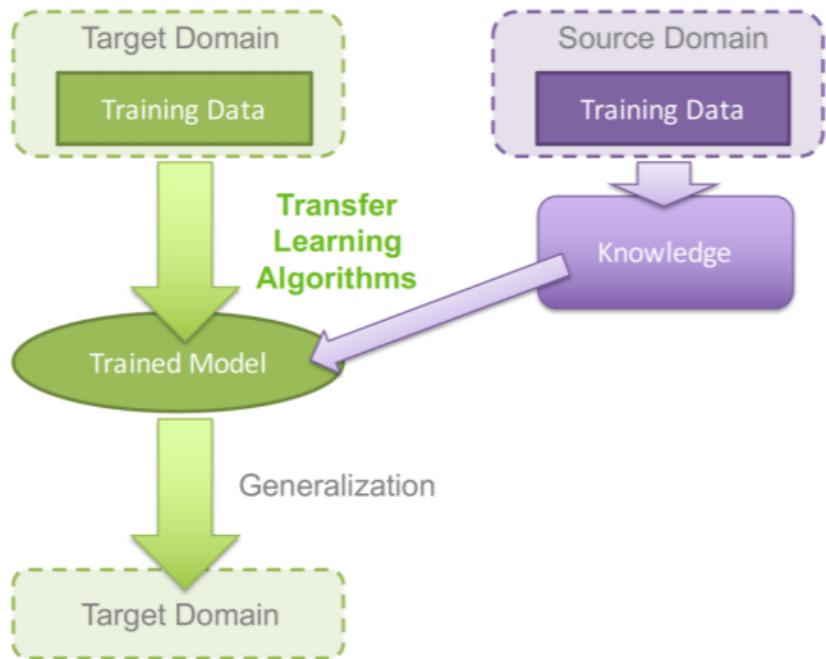
- Proper initialization is an active area of research
- Learning Rate Decay

### Learning Rate and Loss



- SGD, Adagrad, RMSProp, Adam all have learning rate as a hyperparameter. Which one of these learning rates is best to use?
- Answer: All of them! Start with large learning rates and decay over time
- Reduce learning rate at a few fixed points. E.g. for ResNets, multiply  $\alpha$  by 0.1 after epochs 30, 60, and 90.
- Cosine:  $\alpha_t = \frac{1}{2}\alpha_0(1 + \cos(t\pi/T))$
- Linear:  $\alpha_t = \alpha_0(1 - t/T)$
- Transfer Learning

Improve learning  
new tasks by learned  
tasks



- Deep features are fairly transferable, and open-source pre-trained models are now everywhere.

	very similar dataset	very different dataset
very little data	Use linear classifier on top layer	You're in trouble... Try linear classifier from different stages
quite a lot of data	Finetune a few layers	Finetune a large number of layers

## Lecture 14: Deep Learning -- Applications in Low-Level Computer Vision

### Low Level Vision

- Recover the unseen pixels -- image inpainting (修补)

$$\mathbf{y} = \mathbf{Ax} + \mathbf{e} \quad (19)$$

- $\mathbf{x}$ : Signal to be recovered
- $\mathbf{e}$ : Measurement errors
- $\mathbf{A}$ : Sensing matrix
- $\mathbf{y}$ : Measurements



$$(P0) \quad \min_x \underbrace{\|y - Ax\|_2^2}_{\text{Fidelity}}.$$

- What kind of prior knowledge we can have for this ill-condition task
- Deep Learning for Image Restoration and Enhancement
  - Current trend: end-to-end + image-level
  - Future promise: end-to-end + domain-expertise + self-similarity
  - Three examples
    - Image denoising
    - Image deblurring
    - Image super-resolution
  - More "wild" image restoration examples: compression artifact removal, haze/rain/snow removal, underwater/low-light enhancement etc.

$$\mathbf{y} = \mathbf{Ax} + \mathbf{e}$$

$$(P0) \quad \min_x \underbrace{\|y - Ax\|_2^2}_{\text{Fidelity}} + \underbrace{\lambda \Phi(x)}_{\text{Signal Regularizer}}$$

- **x:** Signal to be recovered
- **e:** Measurement errors
- **A:** Sensing matrix
- **y:** Measurements
- **$\Phi(x)$ :** Signal regularizer
  - Sparse coding - dictionary learning
  - Collaborative filtering - BM3D
  - Deep learning - convolutional neural networks

# Image Denoising

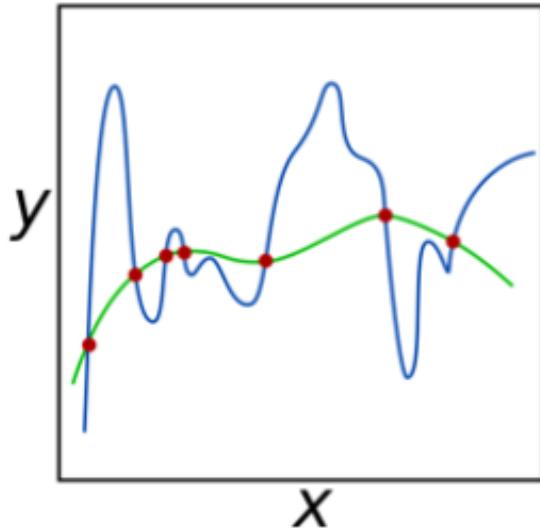
- Image Denoising
  - Simplest low-level vision problem
  - Noisy measurement:  $\mathbf{y} = \mathbf{x} + \mathbf{e}$



- Estimate the clean image:  $\hat{\mathbf{x}} = f(\mathbf{y})$



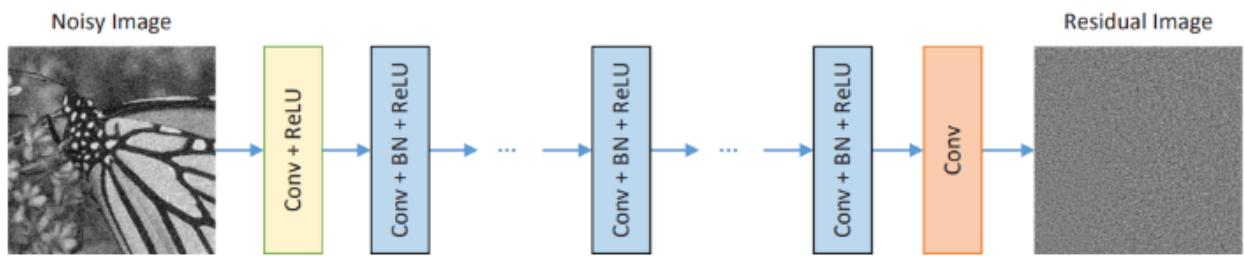
- Conventional(常规) Methods
  - Collaborative (协同) Filtering
    - Non-local means, BM3D, etc.
    - Refs:
      - A. Buades, B. Coll and J. M. Morel, "A non-local algorithm for image denoising," in CVPR, 2005.
      - K. Dabov, A. Foi, V. Katkovnik, and K. Egiazarian, "Image denoising by sparse 3-D transform-domain collaborative filtering," IEEE TIP, 2007.
  - Piece-wise smoothing
    - Total variation, Tikhonov regularization, etc.



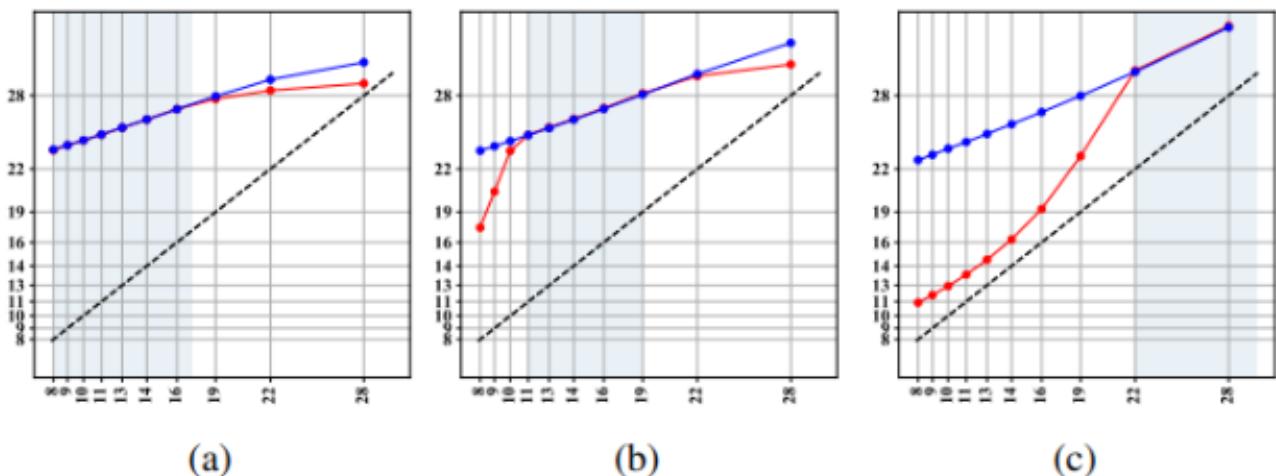
- Sparsity
  - Discrete cosine transform (DCT), Wavelets, etc.
  - Dictionary learning: OMP, Lasso, K-SVD, etc
  - Refs:
    - J. Portilla, V. Strela, M. J. Wainwright, and E. P. Simoncelli, "Image denoising using scale mixtures of Gaussians in the wavelet domain," IEEE TIP, 2003.
    - M. Aharon, M. Elad, and A. Bruckstein, "K-SVD: An algorithm for designing overcomplete dictionaries for sparse representation," IEEE TSP, 2006.
- Image Denoising by Deep Learning
  - Natural idea: train a (convolutional) denoising autoencoder, that regresses clean images from noisy ones
  - It is not easy for deep networks to outperform classical methods such as BM3D
    - BM3D is shown to be better at dealing with self-repeating regular structures
  - How to outperform(在...方面优于...) BM3D using a deep network denoiser?
    - The model capacity is chosen large enough, i.e., enough hidden layers with sufficiently many hidden units
    - The patch (grid) size is chosen large enough, i.e., a patch contains enough information to fit a complicated denoising function that covers the long tail
    - The chosen training set is large enough
- Conventional vs. Deep Learning Methods

	<b>Convenrional</b>	<b>Deep Learning</b>
Model	Shallow model (Equivalently one free layer)	Deep model (Multiple free layers)
Supervise	Unsupervised (1. No training corpus (语料库) needed 2. Data efficient)	Supervised (1. Training corpus needed 2. Data inefficient)
Inverse problem	1. Assumption & understanding of the data 2. Regularizers & structures of the model 3. Flexible	1. Little assumption 2. Almost free model 3. Few work until recently

- [Recent work](#)
- Residual Learning (learn the noisy image)



- MSE is the perfect loss in this setting
- Much easier to optimize
- Ref: K. Zhang, W. Zuo, Y. Chen, D. Meng, and L. Zhang, "Beyond a Gaussian denoiser: Residual learning of deep CNN for imagedenoising,"IEEE TIP, 2017.
- DnCNN with all bias terms removed



- Interpretable (可解释的)
- Robust (强大的)
- Ref: S. Mohan, Z. Kadkhodaie, E. P. Simoncelli, and C. Fernandez-Granda, "Robust and interpretable blind image denoising via bias-free convolutional neural networks,"arXiv.org e-prints, Technical Report 1906.05478, 2019.
- Why are people interested of denoising
  - Simplest inverse problem with  $\mathbf{A} = \mathbf{I}$
  - Ideal test bed to verify the effectiveness of the proposed regularizer

## Image Deblurring

- Image Deblurring
  - Blurred Measurement:  $\mathbf{y} = \mathbf{M} * \mathbf{x}$

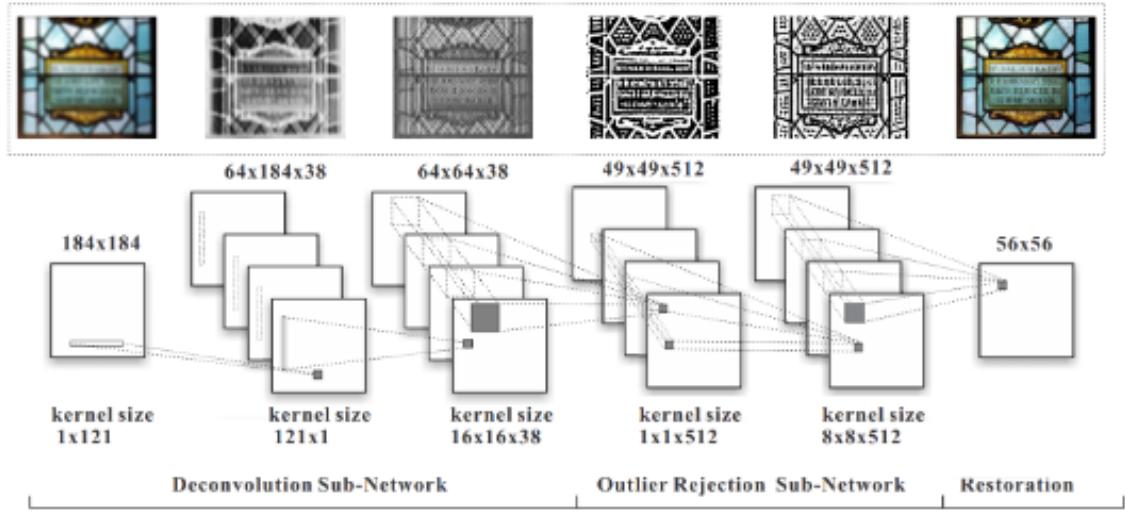
$$\text{Blurry Image} = \text{Sharp Image} * M$$

$$M = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 1 & 1 \\ 0 & 0 & 0 \end{bmatrix}$$

- Where  $M$  is the convolution kernel
- Estimate the stable image:  $\hat{x} = f(y)$

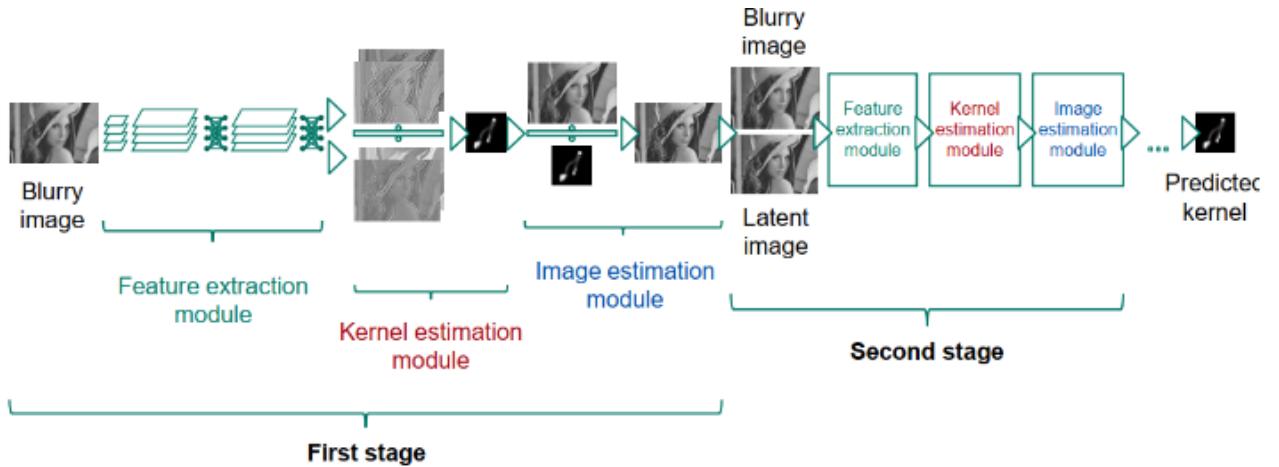


- Non-blind image deblurring
  - Suppose the blurring kernel  $M$  is known
  - $\hat{x} = f(y)$
  - All training data need to have consistent  $M$  as the test data
  - Also difficult because of the compression loss
- Blind image deblurring - a much more challenging problem
  - Estimate both the image, and the blurring kernel
  - $\{\hat{x}, \hat{M}\} = f(y)$
- Image Deblurring by Deep Learning
  - Key technical features:
    - Treat deblurring as a deconvolution task, which can be approximated by a convolution network with larger filter sizes
    - Concatenation of deconvolution CNN module with another denoising CNN module to suppress artifacts and reject outliers (压制伪影并拒绝异常值)



- Ref: L. Xu, J. Ren, C. Liu, and J. Jia, "Deep convolutional neural network for image deconvolution," inNIPS, 2014.

- Blind Image Deblurring



- Key technical features:
  - Iteratively estimate both the blurring kernel and the underlying image
  - Separate networks to capture the image property and the kernel information, respectively
- Ref: C. J. Schuler, M. Hirsch, S. Harmeling, and B. Scholkopf, "Learning to deblur," IEEE TPAMI, 2016.

## Image Super-Resolution

- Image Super-Resolution
  - Low-resolution measurement:  $\mathbf{y} = \mathbf{D} \odot (\mathbf{M} * \mathbf{x})$

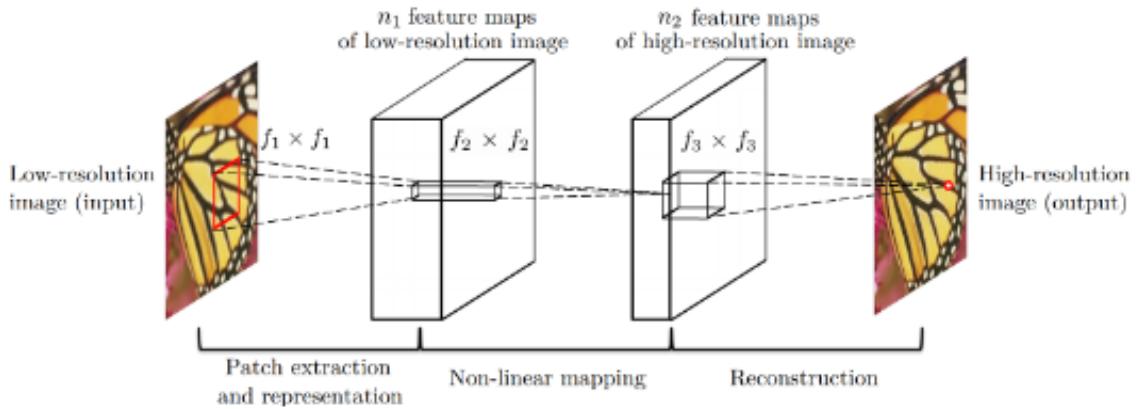
$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \odot \begin{bmatrix} b_{11} & b_{12} & b_{13} \\ b_{21} & b_{22} & b_{23} \\ b_{31} & b_{32} & b_{33} \end{bmatrix} = \begin{bmatrix} a_{11}b_{11} & a_{12}b_{12} & a_{13}b_{13} \\ a_{21}b_{21} & a_{22}b_{22} & a_{23}b_{23} \\ a_{31}b_{31} & a_{32}b_{32} & a_{33}b_{33} \end{bmatrix}.$$



- Estimate the high-resolution image:  $\hat{\mathbf{x}} = f(\mathbf{y})$



- Image Super-Resolution by Deep Learning



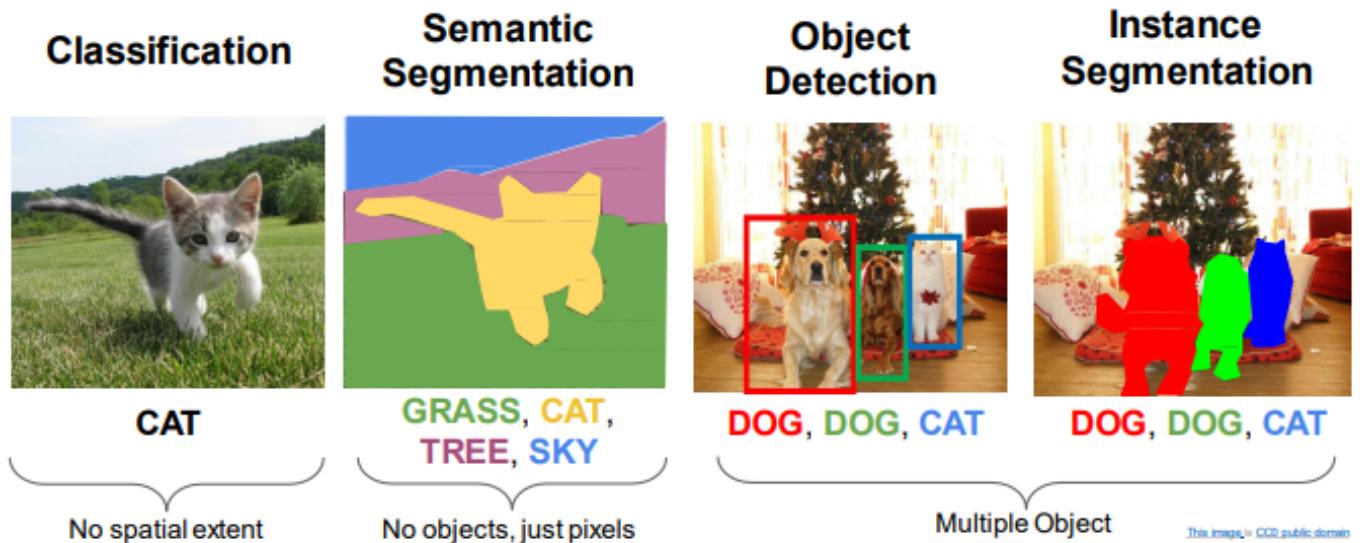
- Key technical features:
  - Learn an end-to-end mapping from low to high-resolution images as a deep CNN
  - Closely mimic the traditional SR pipeline: LR feature extraction => coupled LR-HR feature space mapping => HR image reconstruction
- Ref: C. Dong, C. C. Loy, K. He, and X. Tang, "Image super-resolution using deep convolutional networks," IEEE TPAMI, 2015.

## New Trends

- New topic: dehazing, deraining, low-light enhancement, etc.
- New goal: Human perception v.s. machine perception (知觉)
- New setting: From supervised to unsupervised training (no "GT") or relying on "synthetic pairs"
- New domain: Medical images, infrared images, remote sensing images, etc.
- New concern: "All-in-one" adaptivity, efficient implementation, etc.

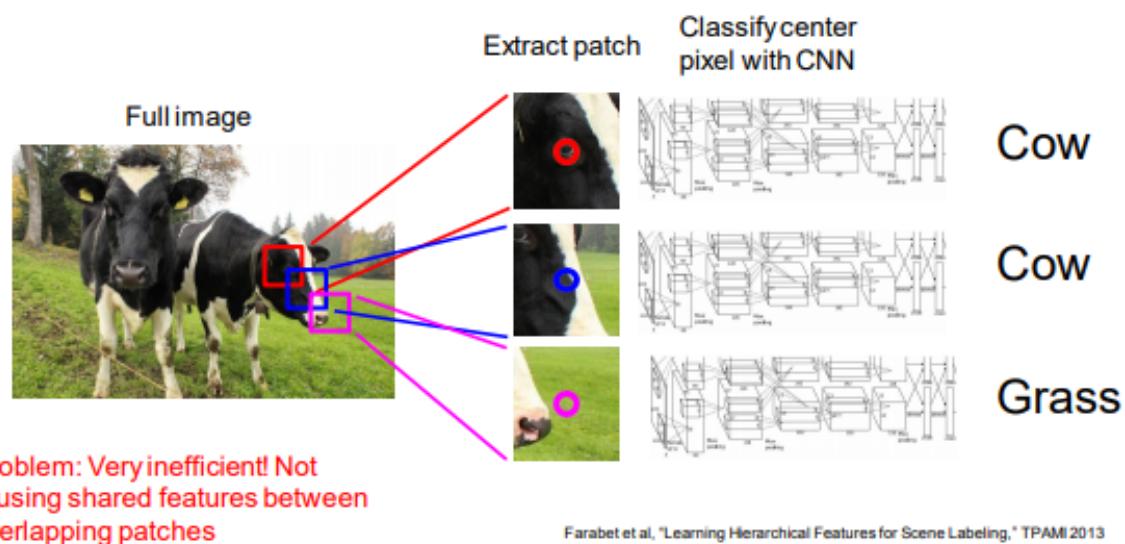
# Lecture 15: Deep Learning -- Applications in High-Level Computer Vision

## Tasks



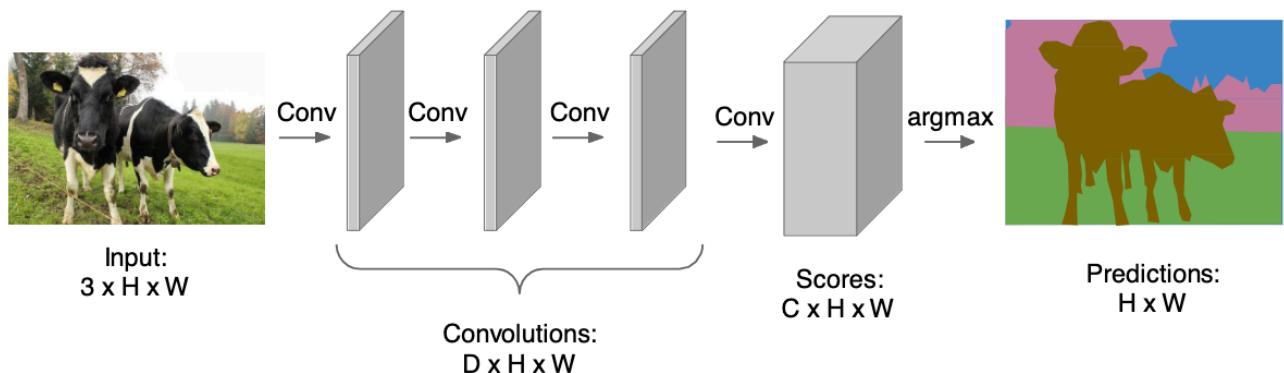
## Semantic Segmentation

- Task:
  - Label each pixel in the image with a category label
  - Don't differentiate instances, only care about pixels
- Sliding Window



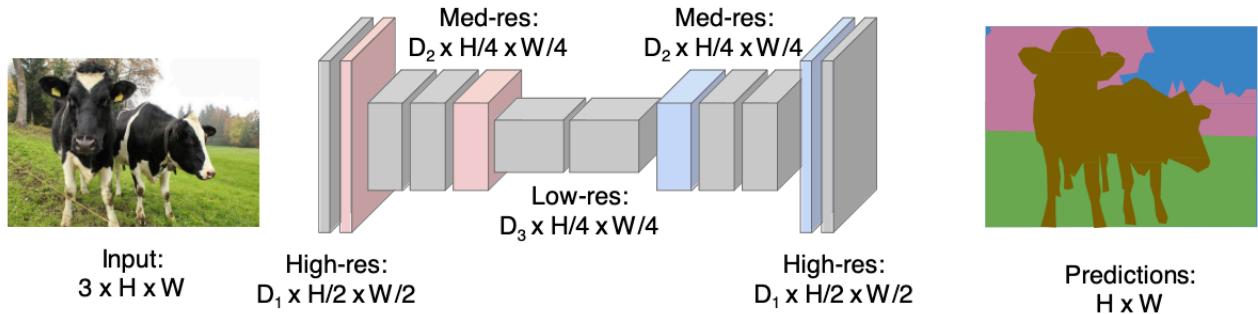
- Fully Convolution

Design a network as a bunch of convolutional layers to make predictions for pixels all at once!



- Problem: Convolutions at original image resolution will be very expensive
- Fully Convolution with downsampling and upsampling

Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!



- $CNN + CNN^{-1}$
- Downsampling: Pooling, strided convolution
- Upsampling: **Unpooling or striped transpose convolution**
- Refs: Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR2015 Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV2015
- Upsampling
  - Unpooling

**Nearest Neighbor**

1	2
1	2
3	4

Input:  $2 \times 2$

1	1	2	2
1	1	2	2
3	3	4	4
3	3	4	4

Output:  $4 \times 4$

**"Bed of Nails"**

1	2
0	0
3	4

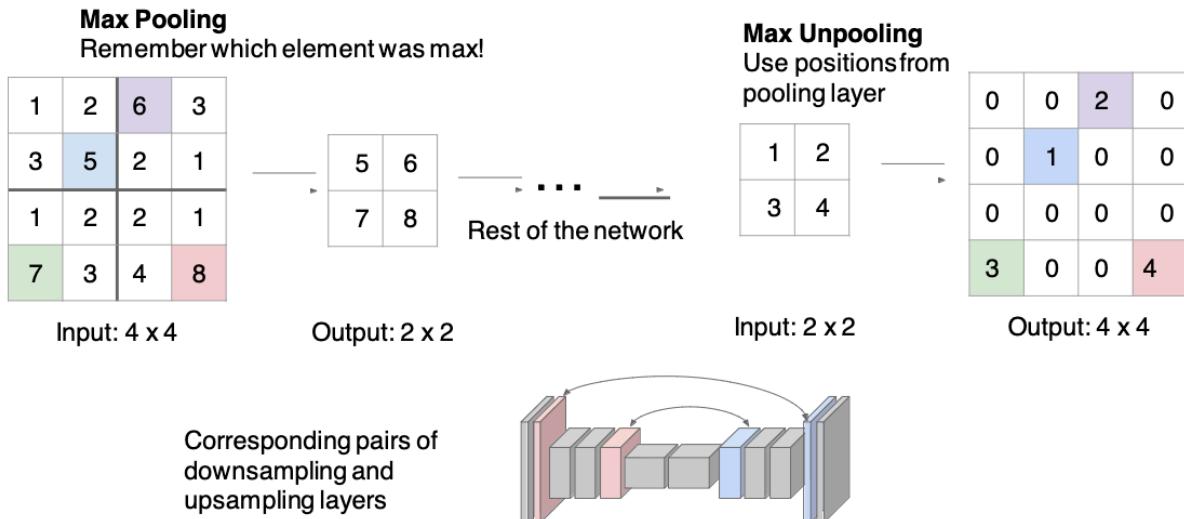
Input:  $2 \times 2$

1	0	2	0
0	0	0	0
3	0	4	0
0	0	0	0

Output:  $4 \times 4$

- Max Unpooling

- Remember the position of copying the value in the max pooling layer, and just copy back to the position and fill zero



## Object Detection

### Object Detection: Impact of Deep Learning

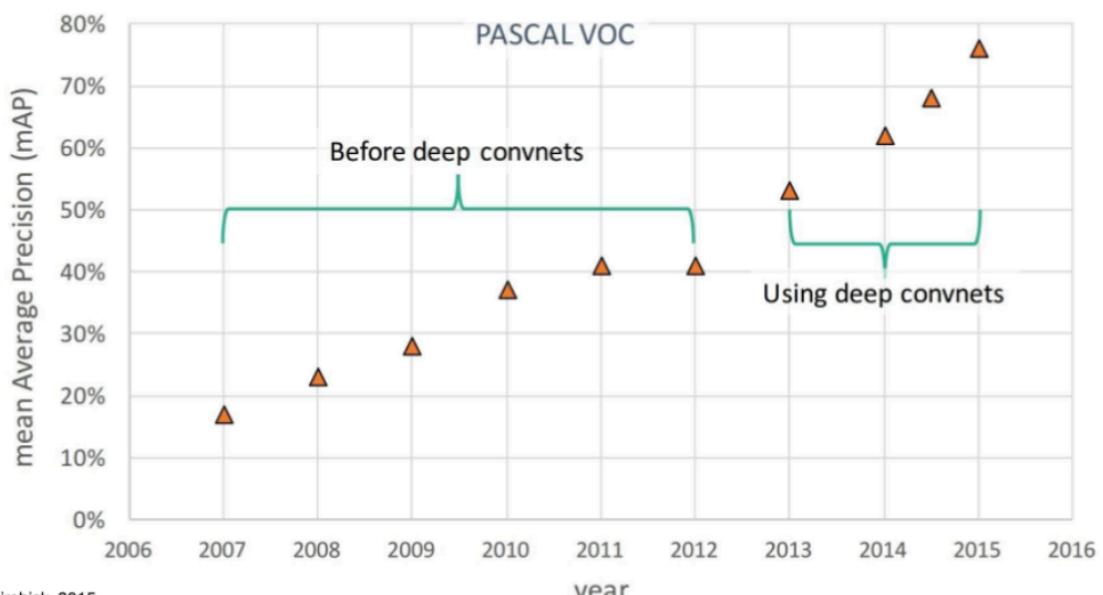
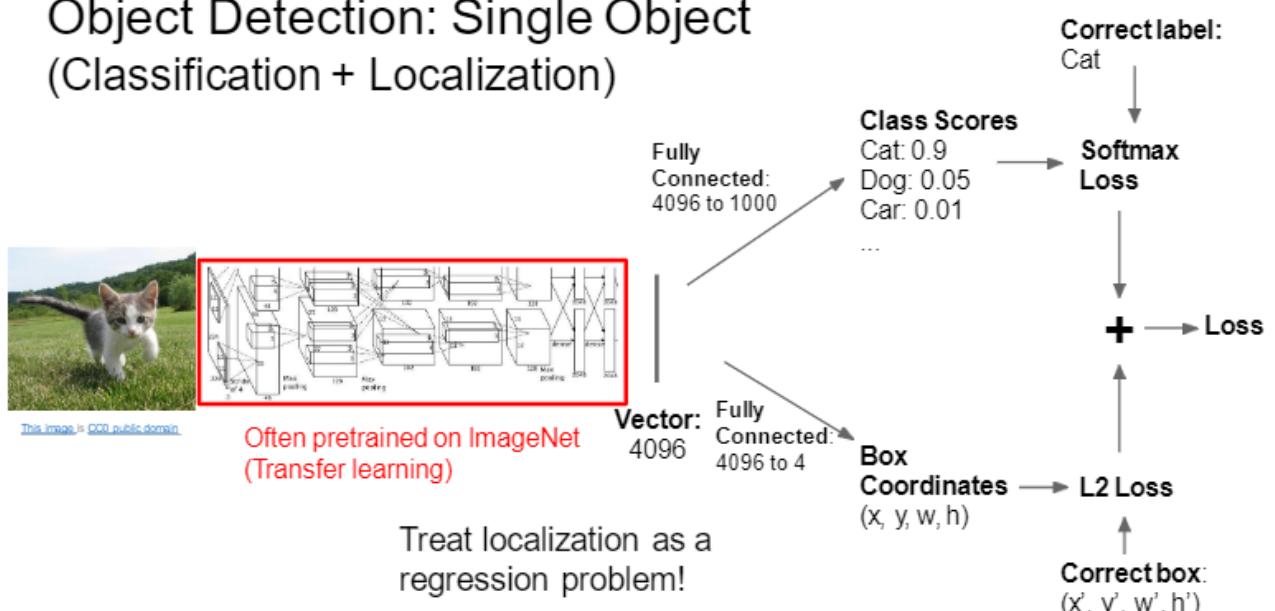


Figure copyright Ross Girshick, 2015.  
Reproduced with permission.

- Single Object Detection => Classification + Localization
  - Use multitask loss

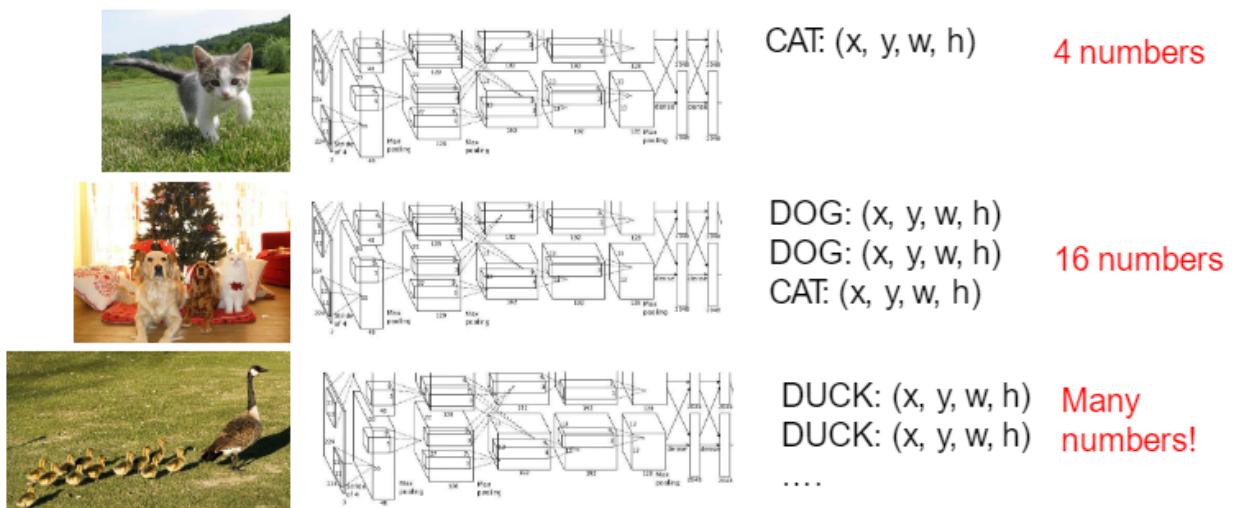
# Object Detection: Single Object (Classification + Localization)



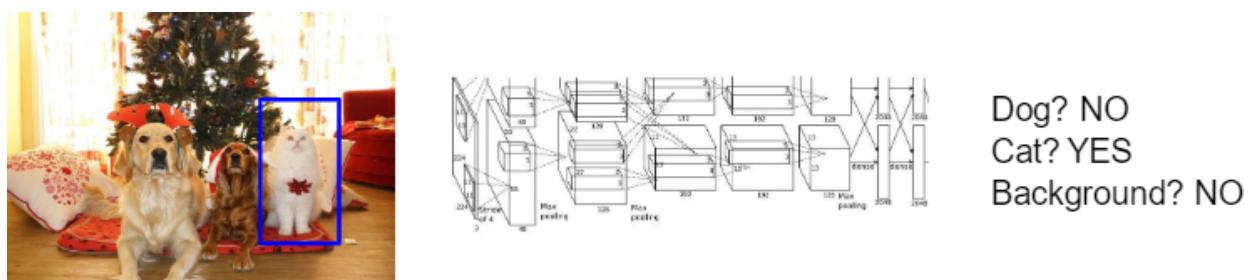
- Multiple Objects Detection => Expensive calculation

# Object Detection: Multiple Objects

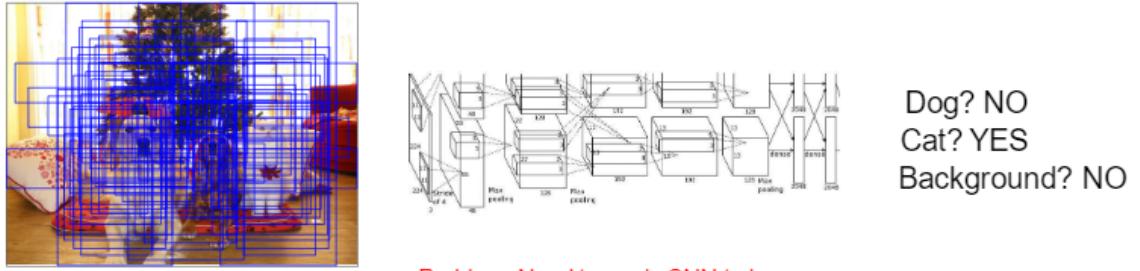
Each image needs a different number of outputs!



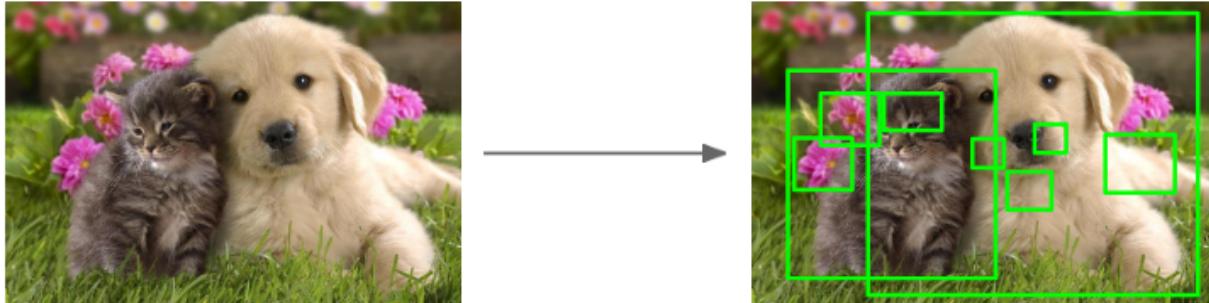
- Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



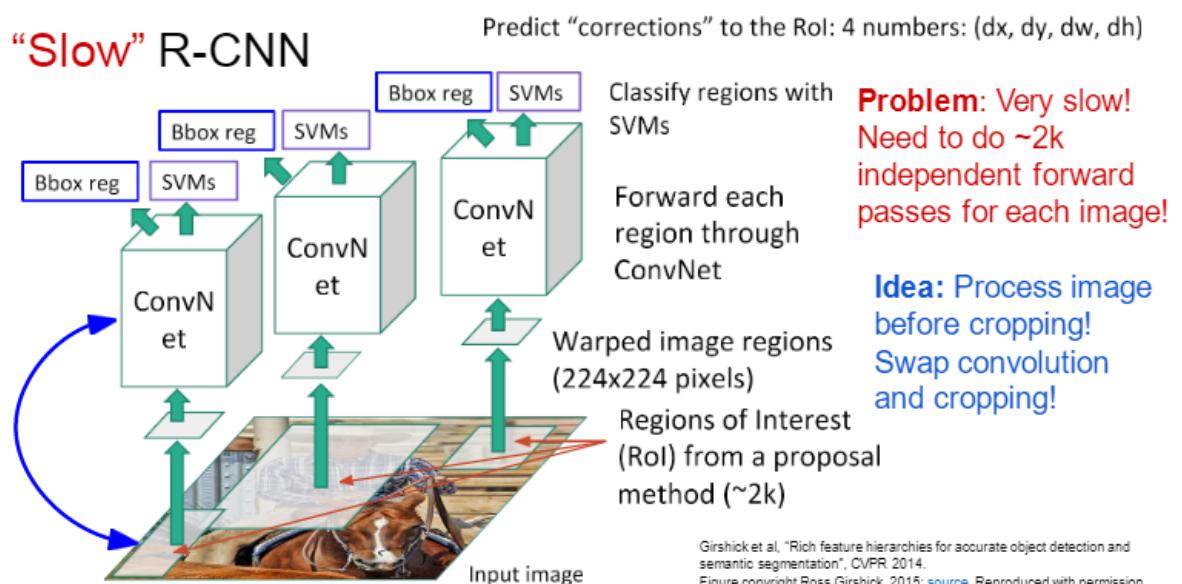
- Problem: Need to apply CNN to huge number of locations, scales, and aspect ratios, very computationally expensive!



- Region Proposals: Selective Search

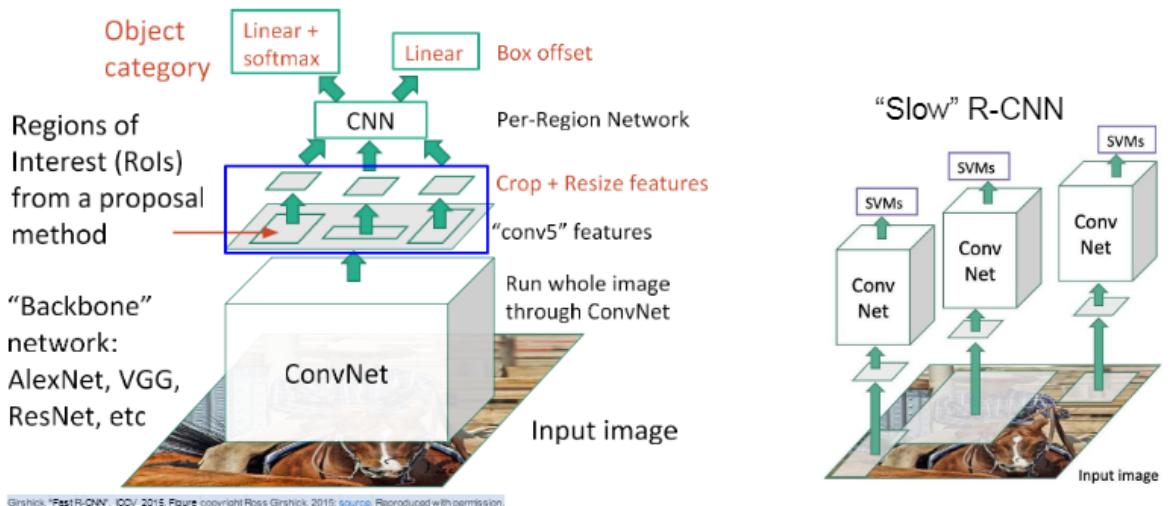


- Find "blobby" that are likely to contain objects
- Relatively fast to run; e.g. Selective Search gives 2000 region proposals in a few seconds on CPU
- Ref: Alexe et al, "Measuring the objectness of image windows", TPAMI 2012 Uijlings et al, "Selective Search for Object Recognition", IJCV2013 Cheng et al, "BING: Binarized normed gradients for objectness estimation at 300fps", CVPR 2014 Zitnick and Dollar, "Edge boxes: Locating object proposals from edges", ECCV2014
- R-CNN
  - "Slow" R-CNN



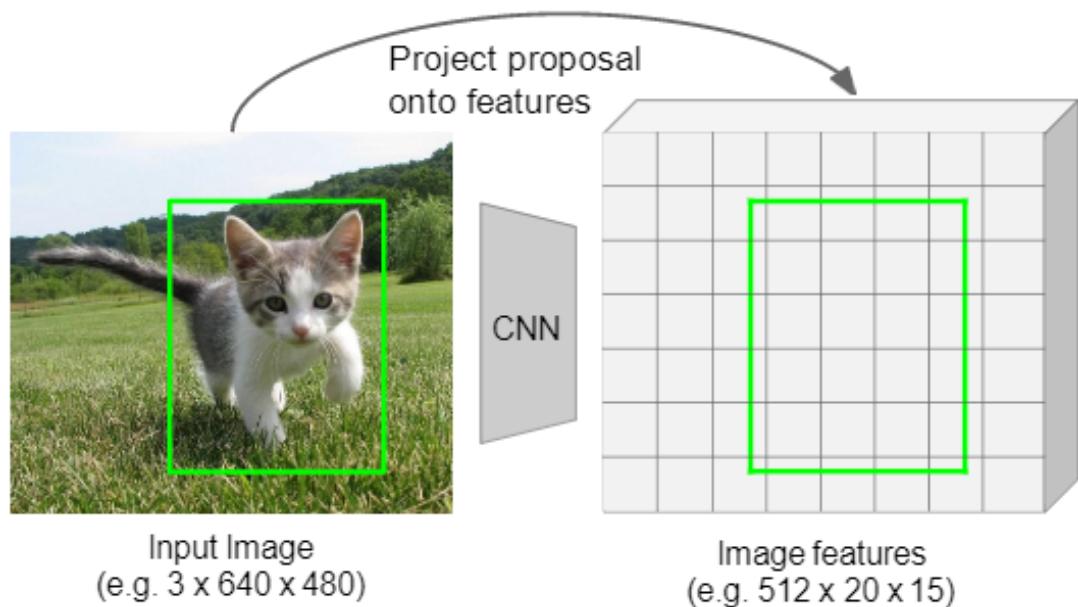
- Ref: Girshick et al, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR2014. Figure copyright Ross Girshick, 2015; source. Reproduced with permission.
- Fast R-CNN

# Fast R-CNN

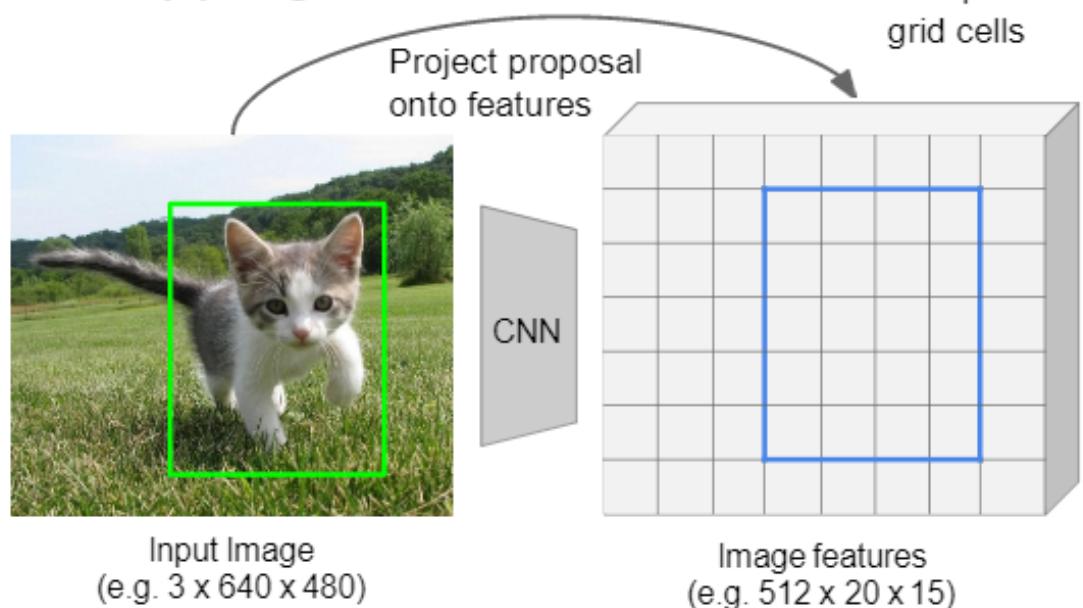


- Feature cropping methods

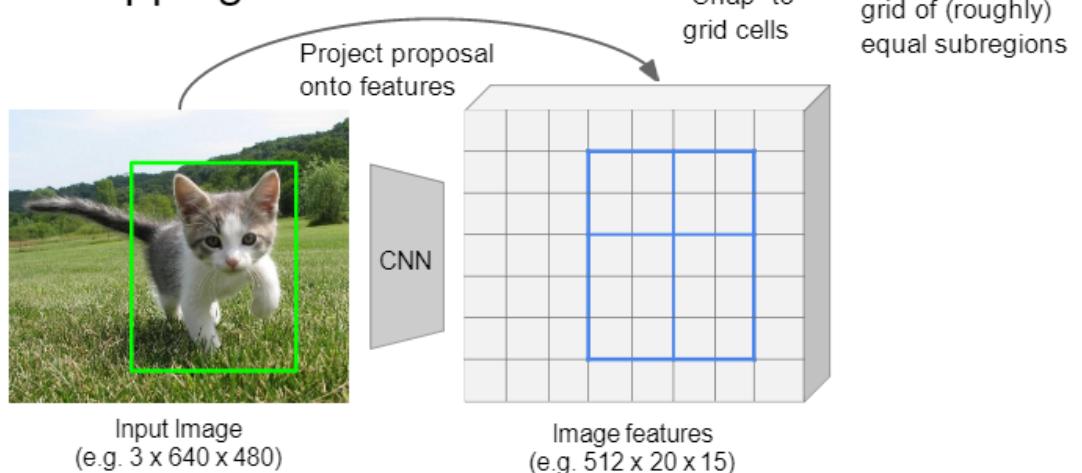
- **Rol Pool**



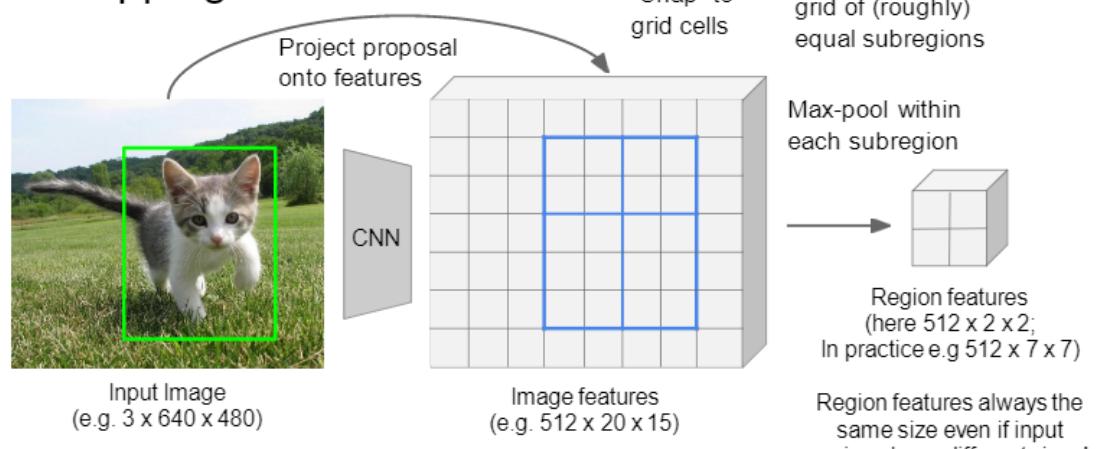
# Cropping Features: RoI Pool



# Cropping Features: RoI Pool

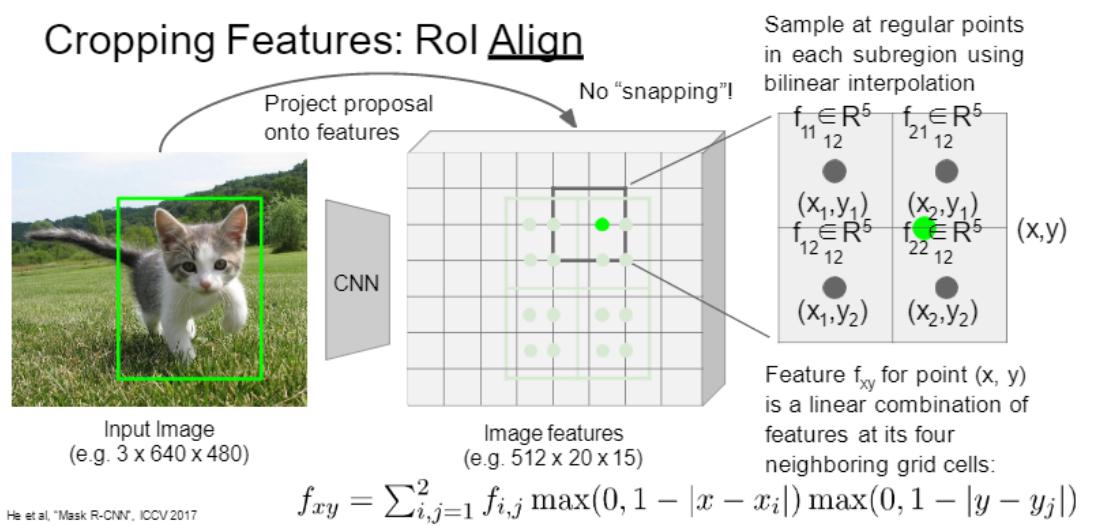


# Cropping Features: RoI Pool



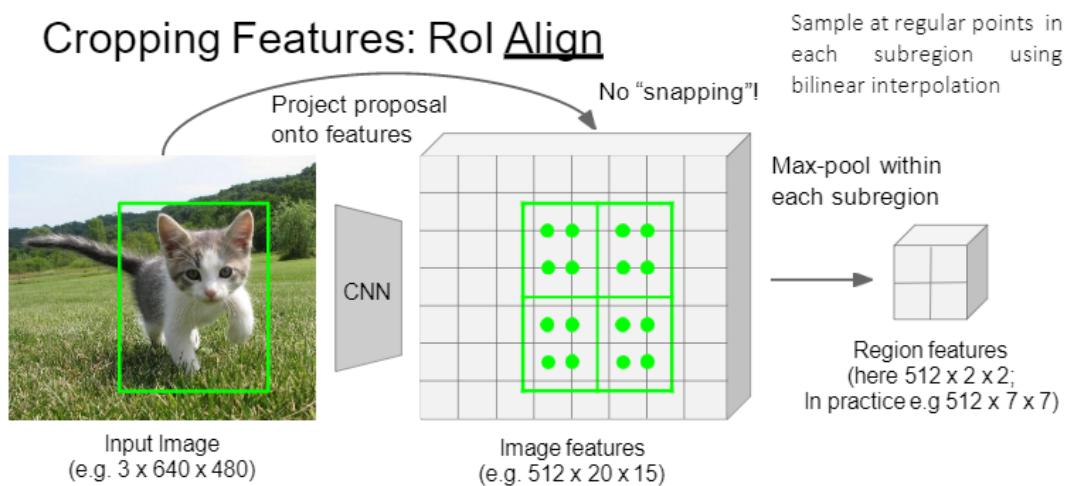
- Problem: Region features slightly misaligned (due to snapping)
- RoI Align

## Cropping Features: RoI Align



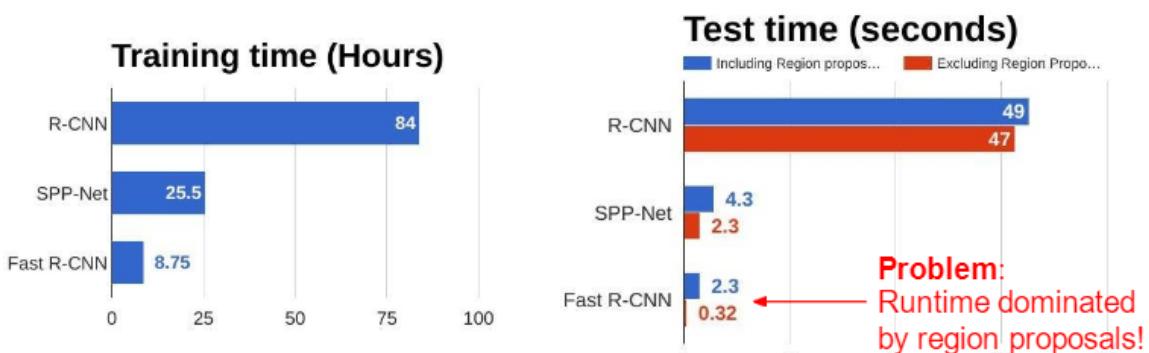
- If the two-axis-distance is less than 1, add the weighted score

## Cropping Features: RoI Align

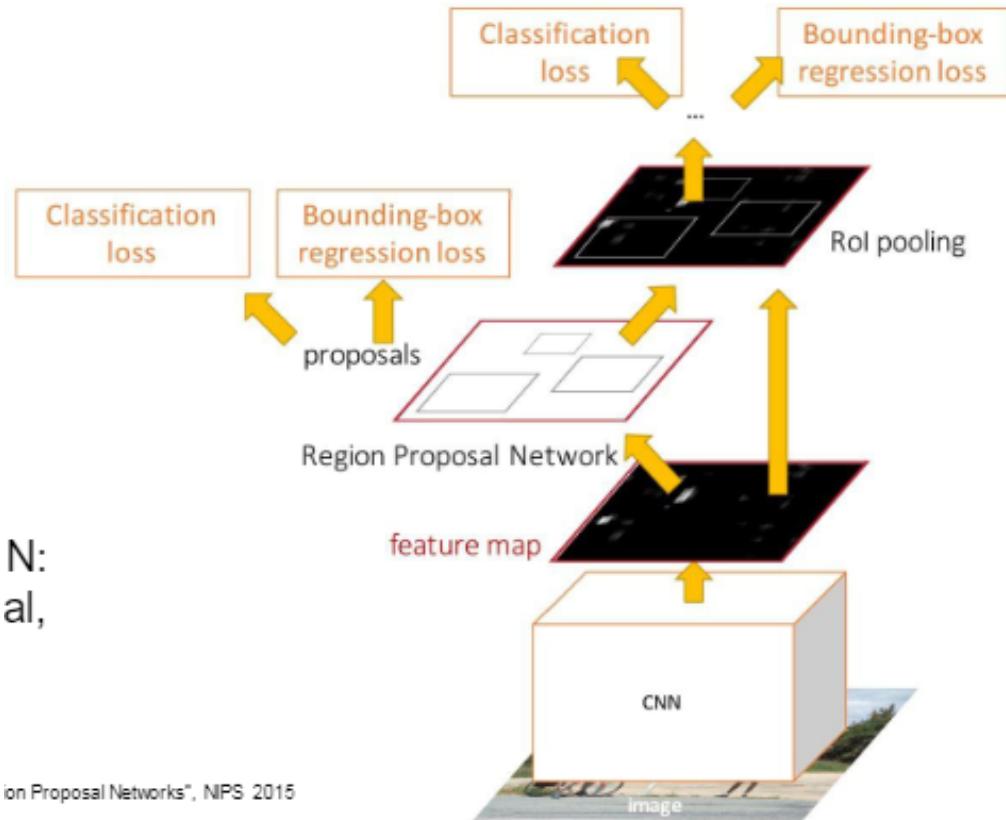


- Ref: Girshick, "Fast R-CNN", ICCV 2015. Figure copyright Ross Girshick, 2015; source. Reproduced with permission.
- R-CNN vs. Fast R-CNN

## R-CNN vs Fast R-CNN



- Faster R-CNN
  - Idea: Make CNN do proposals
  - Method: Insert Region Proposal Network (RPN) to predict proposals from features. Otherwise same as Fast R-CNN: Crop features for each proposal, classify each one. Jointly train 4 losses:
    - RPN classify object / not object
    - RPN regress box coordinates
    - Final classification score (object classes)
    - Final box coordinates

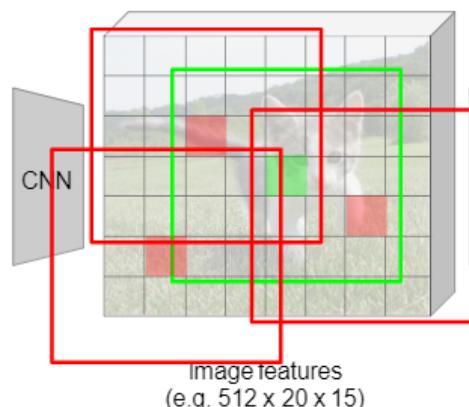


- RPN

## Region Proposal Network



Input Image  
(e.g. 3 x 640 x 480)



Imagine an **anchor box** of fixed size at each point in the feature map

Anchor is an object?  
1 x 20 x 15

At each point, predict whether the corresponding anchor contains an object (per-pixel logistic regression)

# Region Proposal Network



Input Image  
(e.g. 3 x 640 x 480)

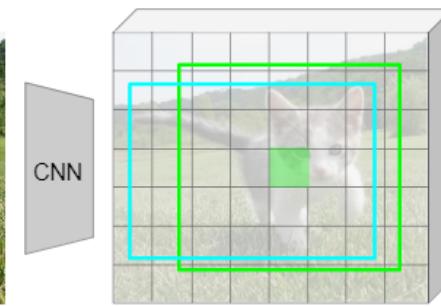


Image features  
(e.g. 512 x 20 x 15)

Imagine an **anchor box**  
of fixed size at each  
point in the feature map

Conv → Anchor is an object?  
1 x 20 x 15  
Conv → Box transforms  
4 x 20 x 15

For positive boxes, also predict  
a transformation from the  
anchor to the ground-truth box  
(regress 4 numbers per pixel)

# Region Proposal Network



Input Image  
(e.g. 3 x 640 x 480)

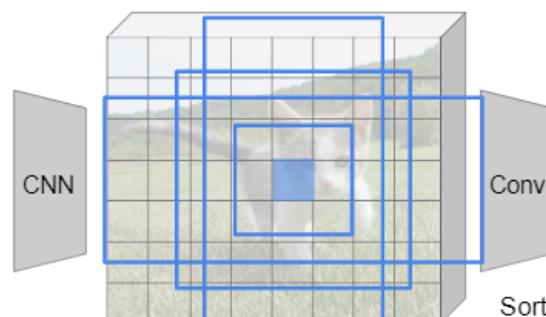


Image features  
(e.g. 512 x 20 x 15)

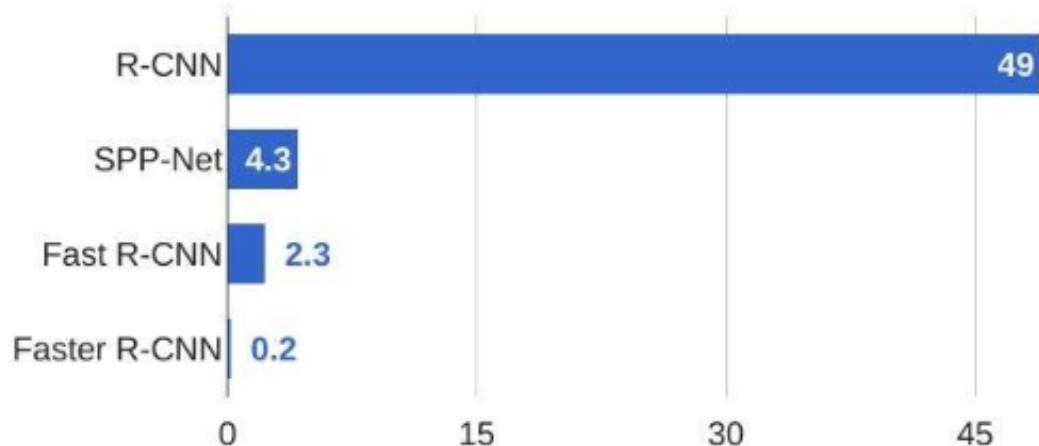
In practice use K different  
anchor boxes of different  
size / scale at each point

Conv → Anchor is an object?  
 $K \times 20 \times 15$   
Conv → Box transforms  
 $4K \times 20 \times 15$

Sort the  $K \times 20 \times 15$  boxes by  
their “object” score, take top  
~300 as our proposals

- Performance

## R-CNN Test-Time Speed



- Glossing over many details

- Ignore overlapping proposals with non-max suppression
- How to determine whether a proposal is positive or negative?
- How many positives/negatives to send to second stage?

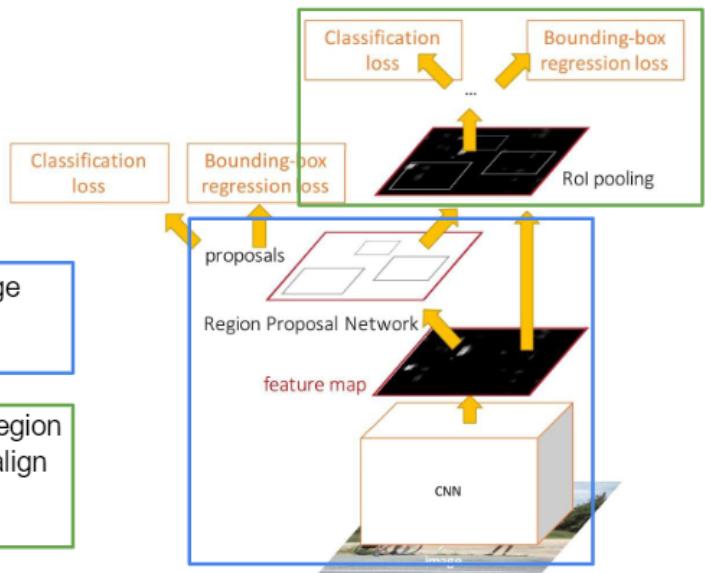
## Faster R-CNN:

Make CNN do proposals!

Faster R-CNN is a  
**Two-stage object detector**

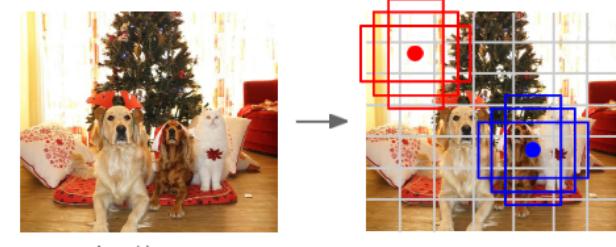
- First stage:** Run once per image
- Backbone network
  - Region proposal network

- Second stage:** Run once per region
- Crop features: RoI pool / align
  - Predict object class
  - Prediction bbox offset



- Do we really need second stage?

## Single-Stage Object Detectors: YOLO / SSD / RetinaNet



Redmon et al., "You Only Look Once: Unified, Real-Time Object Detection", CVPR 2016  
Liu et al., "SSD: Single-Shot MultiBox Detector", ECCV 2016  
Lin et al., "Focal Loss for Dense Object Detection", ICCV 2017

- Within each grid cell:
- Regress from each of the  $B$  base boxes to a final box with 5 numbers:  $(dx, dy, dh, dw, \text{confidence})$
  - Predict scores for each of  $C$  classes (including background as a class)
  - Looks a lot like RPN, but category-specific!
- Output:  
 $7 \times 7 \times (5 * B + C)$

- Ref: Redmon et al, "You Only Look Once: Unified, Real-Time Object Detection", CVPR2016 Liu et al, "SSD: Single-Shot MultiBox Detector", ECCV 2016 Lin et al, "Focal Loss for Dense Object Detection", ICCV2017
- How to parameterize bounding box regression?
- Ref: Ren et al, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015 Figure copyright 2015, Ross Girshick; reproduced with permission
- Comparation

# Object Detection: Lots of variables ...

## Backbone Network

VGG16  
ResNet-101  
Inception V2  
Inception V3  
Inception  
ResNet  
MobileNet

## “Meta-Architecture”

Two-stage: Faster R-CNN  
Single-stage: YOLO / SSD  
Hybrid: R-FCN

## Image Size

## # Region Proposals

...

## Takeaways

Faster R-CNN is slower but more accurate

SSD is much faster but not as accurate

Bigger / Deeper backbones work better

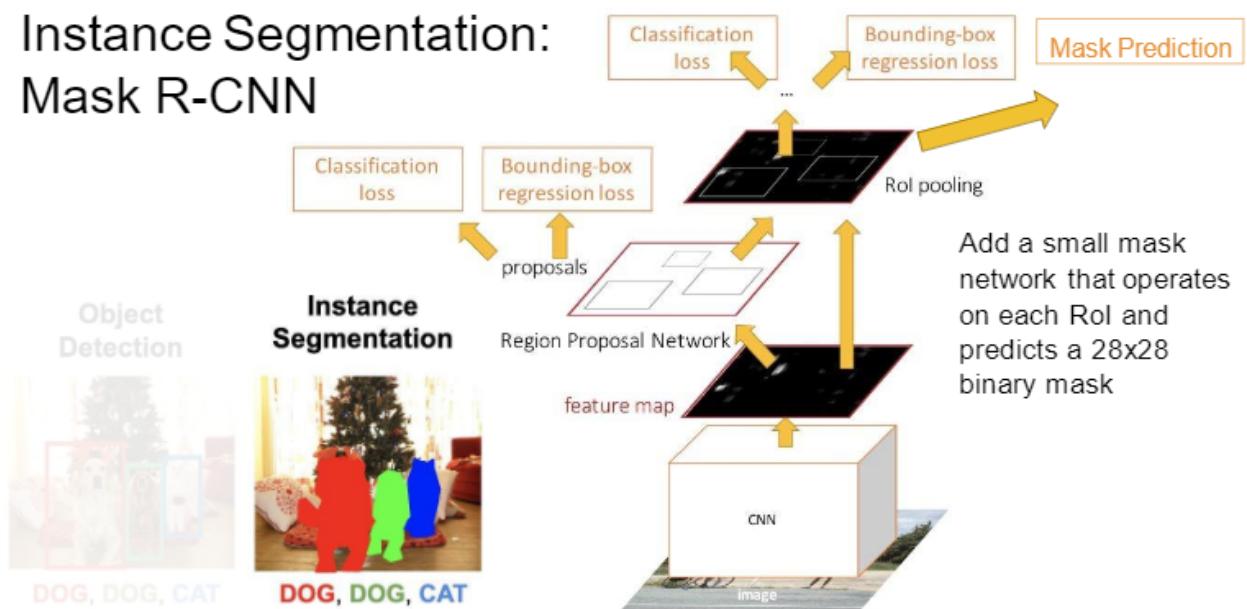
Huang et al, “Speed/accuracy trade-offs for modern convolutional object detectors”, CVPR 2017  
Zou et al, “Object Detection in 20 Years: A Survey”, arXiv 2019 (today!)

R-FCN: Dai et al, “R-FCN: Object Detection via Region-based Fully Convolutional Networks”, NIPS 2016  
Inception-V2: Ioffe and Szegedy, “Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift”, ICML 2015  
Inception V3: Szegedy et al, “Rethinking the Inception Architecture for Computer Vision”, arXiv 2016  
Inception ResNet Szegedy et al, “Inception-ResNet and the Impact of Residual Connections on Learning”, arXiv 2016  
MobileNet: Howard et al, “Efficient Convolutional Neural Networks for Mobile Vision Applications”, arXiv 2017

# Instance Segmentation

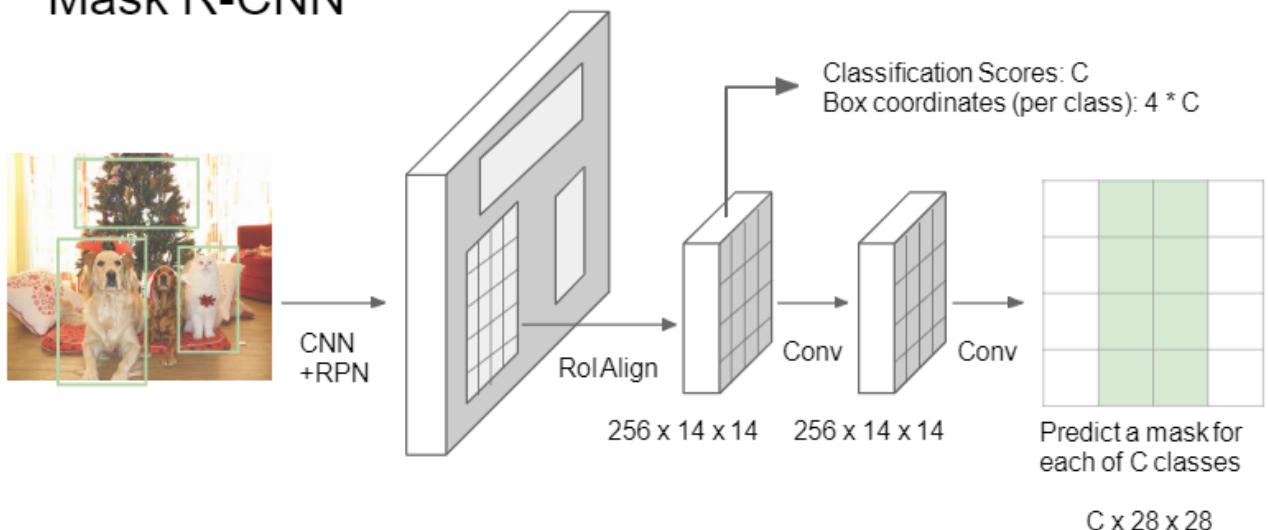
- Mask R-CNN

## Instance Segmentation: Mask R-CNN



He et al, “Mask R-CNN”, ICCV 2017

# Mask R-CNN



He et al, "Mask R-CNN", arXiv 2017