

# L2 - Supervised ML

October 20, 2024 1:05 AM

- mapping a given input to a target
- i.e. if given a dataset where each data instance is  $x_i$ , and the target is  $y$ , then the goal of supervised learning is to find  $f$  s.t.
  - o  $f(x_1, x_2, \dots, x_p) \rightarrow y$

f	model
$x_1, \dots, x_p$	dataset
y	target

- We build the predictive model using a dataset and an algorithm
- Learner: the statistical model that we are learning for the task
- Observation: the single data instance of the model
- Training set: the set of estimation samples used for computing the model
- Test set: the out-of-sample observation that the learner has not seen before
- Algorithm: the estimation method used to compute the model
- Features: the vector representation of a single data instance used in the algorithm
- Overfitting: models following errors (noise) too closely
  - o happens when model is too flexible/complex

## Prediction and inference

- two reasons to estimate a model ( $f$ )
- Prediction
  - o when a set of  $X$  is available but the output  $Y$  cannot be easily obtained
  - o the error term averages to zero, therefore we can predict  $Y$  with  $Y = g(X)$ , where  $g$  is an estimation of the actual function  $f$
  - o reducible error
    - can be corrected by choosing the most appropriate technique to estimate  $f$
  - o irreducible error
    - variability associated with the error term in  $X$  (epsilon)
- Inference
  - o understand the association between  $X$  and  $Y$
  - o goal is to estimate  $f$  but not necessarily predict  $y$
  - o  $f$  is not a black box, we want to know its exact form
  - o Which predictors are associated with the response?
  - o What is the relationship between the response and each predictor?
  - o Can the relationship between  $Y$  and each predictor be summarized using a linear equation, or is it more complicated?
  - o i.e. what tv do people purchase based on price, store location, discount, competition price

## Parametric models

- a learning model that summarizes data with a set of parameters of fixed size
- line of best fit
- ex:
  - o linear discriminant analysis
  - o perceptron
  - o naïve bayes
  - o simple neural networks

- of the form  $f(X) = B_0 + B_1X_1 + B_2X_2 \dots B_pX_p$  assuming  $f$  is linear in  $X$
- makes it so that we only need to estimate the coefficients, instead of a  $p$ -dimensional function  $f(X)$ 
  - o can estimate the  $B$ 's using OLS
- Pros:
  - o speed: very fast to learn
  - o less data needed: no large datasets necessary and can work well with imperfect data
- cons:
  - o constrained to their very specific form
    - used only for simpler problems
  - o not very accurate
    - unlikely to match the underlying mapping function
    - usually doesn't match the true form of  $f$

#### Non-parametric models

- do not make a lot of assumptions about the data
- examples:
  - o k-nearest neighbors
  - o decision trees
  - o SVM
- pros:
  - o flexibility: fit many functional forms
  - o power: no assumptions/weak assumptions about the underlying function
  - o performance
- cons:
  - o more data needed to estimate mapping function
  - o slower to train, far more parameters
  - o overfitting: high risk of overfitting
  - o explainability: harder to explain why specific predictions are made

#### Accuracy vs Interpretability

- linear regression is inflexible because it can only generate linear functions
- Why use a more restrictive approach vs a flexible one?
- For inference, restrictive models are much more easily interpretable
  - o easy to show relationships between  $X_i$  and  $Y$
- Very flexible methods can lead to very complicated estimates of  $f$
- For prediction not inference, we don't care about interpretability
  - o however, sometimes we get better results using inflexible methods
    - because of overfitting

#### Linear regression

$$y' = \beta_0 + \beta_1 * x_1 + \beta_2 * x_2 + \dots + \beta_p * x_p$$

- $y_1' = \beta_0 + \beta_1 * x_{11} + \beta_2 * x_{12} + \dots + \beta_n * x_{1p}$   
 $y_2' = \beta_0 + \beta_1 * x_{21} + \beta_2 * x_{22} + \dots + \beta_n * x_{2p}$   
 $y_n' = \beta_0 + \beta_1 * x_{n1} + \beta_2 * x_{n2} + \dots + \beta_n * x_{np}$

$B_i$	feature weights
$x_i$	feature value
$n$	number of features

- $$X = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{pmatrix}$$

$p$  = number of features  
 $n$  = number of observations

- $$Y' = h_{\beta}(x) = \beta \cdot X$$

$h_{\beta}$  is the hypothesis function       $\beta$  is model parameters

## Ordinary least squares

- find  $\beta$  such that it minimizes  $\sum(d^2)$  where  $d$  is the distance from each datapoint to the line of best fit

- | Input data  | Target  | Objective                                       |
|---|---|---|
| $X = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{pmatrix}$ | $y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$ | $MSE = \frac{1}{n} \sum_{i=1}^n (y_i - y'_i)^2$ |
| Loss function (K) = Mean Squared Error  |   |   |

- closed form of loss function

- Closed Form Equation
  - $$\beta' = (X^T X)^{-1} X^T y$$

## Other loss functions

- RMSE

- $$\sqrt{\frac{1}{n} \sum_{i=1}^n (y'_i - y_i)^2}$$

- MAE

- $$\frac{1}{n} \sum_{i=1}^n |y'_i - y_i|$$

# L3 - gradient descent

October 20, 2024 5:43 PM

OLS

- Loss Function (K) =  $\frac{1}{n} \sum_{i=1}^n (y_i - y_i')^2$

- Find beta such that it minimizes the loss
  1. Find the derivative of K wrt B
  2. Make it equal 0

- $\nabla_{\beta} K(\beta) = 0 \Rightarrow \dots \Rightarrow \nabla_{\beta} \frac{1}{2} (X\beta - y)^T (X\beta - y) = 0 \Rightarrow \dots \Rightarrow \beta = (X^T X)^{-1} X^T y$

3. Predicted value  $B' = (X^T X)^{-1} X^T y$

- What do we do if there are too many variables?

Gradient descent

- continuously check for local minimum and head in the direction of the negative gradient
- To find the minimum value of the loss function K
  1. obtain the gradient of K
  2. follow the gradient downhill
    - i. Start with any arbitrary value of Beta
    - ii. find the steepest descent =  $\text{Grad}(B) * K$
    - iii. compute the new value  $B_{\text{new}} = B - \alpha * \text{grad}(B) * K$
    - iv. go to 2(ii) and repeat
  - steepest descent converges when every element of the gradient is zero:  $\nabla_B * K = 0$
  - $\text{grad}(B) * K =$  vector of partial derivatives w.r.t. each beta

$$\nabla_{\beta} f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial \beta_1} \\ \frac{\partial f(x)}{\partial \beta_2} \\ \vdots \\ \frac{\partial f(x)}{\partial \beta_p} \end{bmatrix}$$

- for linear regression....

- $K = \frac{1}{n} \sum_{i=1}^n (y_i - y_i')^2 = \frac{1}{n} \sum_{i=1}^n (y_i - B * x_i)^2$

- $\nabla_B * K = \nabla_B * \left( \sum_{i=1}^n (y_i - B * x_i)^2 \right) = \frac{2}{n} (B * x - y) * x^T$

# L3-II - polynomial regression

October 20, 2024 6:20 PM

## Polynomial regression

- assume non-linear relationship between variables
- Our linear function:  $y' = B_0 + B_1 * x_1 + \text{error}$
- our non-linear (n-degree) polynomial function

$$y' = \beta_0 + \beta_1 * x_1 + \beta_2 * x_1^2 + \beta_3 * x_1^3 + \dots + \beta_n * x_1^n + \epsilon$$

- different from multiple linear
  - $y' = B_0 + B_1 * x_1 + B_2 * x_2 + \dots + B_p * x_p$ 
    - o powers of  $x_1$  \* each beta in n-degree
    - o  $x_1$  \* one beta,  $x_2$  \* another, ... for multiple linear

- How to learn parameters for poly regression

1. hypothesize a non linear model

$$y' = \beta_0 + \beta_1 * x_1 + \beta_2 * x_1^2 + \beta_3 * x_1^3 + \dots + \beta_n * x_1^n + \epsilon$$

2. select a loss function i.e.  $MSE = (1/n) * \sum((y_i - y_i')^2)$

3. Find beta s.t. it minimizes K

$$\beta' = (X^T X)^{-1} X^T y$$

- beta can be easily estimated using OLS because hypothesis is just a standard linear model with predictors  $x_p, x_p^2, x_p^3 \dots x_p^n$

$$\beta' = (X^T X)^{-1} X^T y$$

$$X = \begin{bmatrix} x_1 & x_1^2 & \dots & x_1^n \\ x_2 & x_2^2 & \dots & x_2^n \\ \vdots & \vdots & \dots & \vdots \\ x_p & x_p^2 & \dots & x_p^n \end{bmatrix} \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_p \end{bmatrix}$$

o

$p$  = number of observations

$n$  = degree of polynomial

- in polynomial functions, flexibility increases as  $n$  increases

# L4 - evaluation

October 20, 2024 6:29 PM

- to evaluate a model we need to quantify how close the predicted response value is to the true response value
- for regression, this difference is called error
  - o should be very small for each observation
- we care about the accuracy on the unseen instances (test set)
- If we had a large number of test observations, we would compute the test MSE

$$\text{Test MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - y_i')^2$$

- When only training data is available: use training MSE
- When training and test data are available: use test MSE
- Overfitting: when small training MSE yields large test MSE
  - o we always expect training MSE to be smaller than the test MSE
  - o overfitting is when a less flexible model would have yielded a smaller test MSE
- Underfitting: when a high training MSE yields a high test MSE
  - o underfitting is when a model is too simple and fails to capture relationships between the i/o
  - o result of a model needing more training time, more input features, less regularization

## Bias variance tradeoff

- every model has an error term
  - o  $y' = \beta_0 + \beta_1 * x_1 + \beta_2 * x_1^2 + \beta_3 * x_1^3 + \dots + \beta_n * x_1^n + \epsilon$
- For a given observation  $x_0$ ,  
**expected test MSE = Variance of  $y'$  + (Bias of  $y'$ )<sup>2</sup> + Variance of  $\epsilon$**
- variance is the amount by which  $y'$  would change if we trained it using a different data set
- bias is the amount of error in  $y'$  when applied to real-world problems
  - o more flexible methods = less bias
  - o i.e. training using linear regression when the relationship between  $x$  and  $y$  is highly nonlinear = high bias
    - no matter how many training observations, it will not be possible to produce an accurate estimate using linear regression
- variance of  $\epsilon$  is irreducible
- **the bias-variance tradeoff**: good test set performance requires low variance as well as low squared bias
  - o underfitting: simple models trained on different samples of the data do not differ much from each other
  - o high variance, overfitting: complex models trained on different samples of the data are very different from each other

## Cross validation

- can be used to estimate the test error associated with a learning method to evaluate its performance
- a good model is one whose test error rate is small
- Validation sets
  - o divide the data into training and validation set
  - o train on the training set, compute performance on the validation set
  - o cons:

- inconsistent test MSE performance - depends on which observations are in the training set
    - bias
  - number of training set observations reduces significantly
- Leave out one (LOOCV)
  - divide into training set and one validation observation
  - train on training set, compute performance on the single observation
  - repeat with every observation
  - pros:
    - less bias
    - no randomness in split
  - cons
    - computation heavy
    - highly variable
- K folds
  - randomly divide the set of observations into k groups (folds) of approximately equal size
  - first fold is treated as a validation set, remaining k-1 folds are used for training
  - tradeoff between loocv and splitting into one training and one validation set
    - loocv has slightly higher variance than k fold, but lower bias
  - typically we choose k = 5, 10 because they have been shown to suffer from neither excessively high bias nor from high variance

#### Error

- test error - average error that results from using a learning method to predict the response on a new observation
  - the use of a particular method is warranted if it results in low test error
- training error - average error resulting from using a learning model on the test set
-

# L5 - classification and logistic regression

October 20, 2024 7:27 PM

- predicting where y is a categorical variable

## Binary classification

- yes or no problems
- linear model eqn:  $y = \beta_0 + \sum(\beta_i * x_i) + \epsilon$ 
  - o right side is continuous unbounded, left side is binary
    - y is 0 or 1
  - o predict the probability that a given instance belongs to 1 class of y
    - $P(y = 1|x) \propto \beta_0 + \sum_{i=1}^p \beta_i * x_i + \epsilon$ 
      - compares 2 inputs and identifies which one leads to a higher probability of y belonging to a given class
        - ◆ we can rank these probabilities
  - o  $\beta$  is a parameter matrix and x represents features, so we can rewrite as
    - $P(y = 1|x) = b + W * X$ 
      - sum of the weighted features (dot product), adding them to a bias term
      - our algorithm needs to learn W and b
      - does not provide probability, still looks like a linear model
  - o squash outputs to (0, 1) using the logistic function (sigmoid function)
    - $\sigma(z) = \frac{1}{1 + e^{-z}}$ 
      - has the property  $1 - \sigma(z) = \sigma(-z)$
  - o Apply logistic function to our linear model
    - $P(y = 1|x) = \sigma(b + W * X)$
  - o simplifies to

$$\log\left(\frac{p(x)}{1 - p(x)}\right) = \beta_0 + \sum_{i=1}^p \beta_i x_i$$

where  $p(x) = b + W * X$

logit function = logistic regression model

- can predict y=1 if probability is high and y=0 else
- set threshold, usually 0.5

## Organizing data

$$\log\left(\frac{p(x)}{1 - p(x)}\right) = \beta_0 + \sum_{i=1}^p \beta_i x_i \Rightarrow y = b + W \cdot X$$

p = number of observations  
n = number of features

-

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_p \end{bmatrix} \quad X = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \dots & \vdots \\ x_{p1} & x_{p2} & \dots & x_{pn} \end{bmatrix} \quad W = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_p \end{bmatrix} \quad b = \begin{bmatrix} b \\ b_2 \\ \vdots \\ b_p \end{bmatrix}$$

$p * 1 \qquad p * n \qquad p * 1 \qquad p * 1$

$$y = X \cdot \text{dot}(W) + b \quad \text{To align our vectors for multiplication}$$

## Logistic regression

- how to train?
- select a cost function: the likelihood function
  - o probability of y = 1 and y = 0 given p(x)
    - $P(x_i, y_i | \beta) = \begin{cases} p(x_i) & , \text{if } y_i = 1 \\ 1 - p(x_i) & , \text{if } y_i = 0 \end{cases}$
- generalize and apply log to both sides gives log-likelihood function for entire dataset
  - $\ell(\beta) = \sum_{i=1}^N \log(P(x_i, y_i | \beta)) = \sum_{i=1}^N \{y_i \log(p(x_i)) + (1 - y_i) \log(1 - p(x_i))\}$
- goal is to learn W and b to maximize the log probability of correct label p(y|x) in the training data
  - $\log(p(y_i|x_i)) = -(y_i \log(p(x_i)) + (1 - y_i) \log(1 - p(x_i)))$

## How to train logistic regression?

1. Cross entropy loss
  - o ensures the probability of the correct answer is maximized and probability of incorrect answer is minimized
  - o the log likelihood loss function



$$\log(p(y_i|x_i)) = y_i \log(p(x_i)) + (1 - y_i) \log(1 - p(x_i))$$

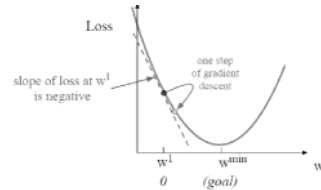
- the loss function needs to be minimized
  - use the cross-entropy loss (also called negative log likelihood)
  - sub  $p(x) = \sigma(b + W * X)$
$$\log(p(y_i|x_i)) = -(y_i \log(\sigma(b + W * X)) + (1 - y_i) \log(1 - \sigma(b + W * X)))$$

- i.e. if  $b = 10.875$  and  $w = -0.171$   
 if the correct label for  $x = 31$  is  $y = 1$ , then CE loss =  $\log(\sigma(b + W * x)) = \log(\sigma(10.875 - 0.171 * 31))$   
 if the correct label for  $x = 31$  is  $y = 0$ , then CE loss =  $-\log(1 - \sigma(b + W * x)) = -\log(1 - \sigma(10.875 - 0.171 * 31))$

## 2. We need an optimization algorithm

- we know the sigmoid function is differentiable
- we know logistic regression is convex
  - it has at most one minimum
  - there is no local minimum, only one global min
    - can't get stuck

ideal for Gradient Descent



- cross entropy loss (L)  

$$L = -(y_i * \log(\sigma(b + w * x)) + (1 - y_i) * \log(1 - \sigma(b + w * x)))$$

for a single observation:

gradient of CE loss

$$\frac{\partial L}{\partial w_i} = [\sigma(b + w * x) - y] * x_i = (y' - y) * x_i$$

$$\frac{\partial L}{\partial b} = [\sigma(b + w * x) - y] = (y' - y)$$

- for many observations

$$L = -\frac{1}{n} \sum_{i=1}^n y_i * \log(\sigma(b + w * x)) + (1 - y_i) * \log(1 - \sigma(b + w * x))$$

gradient of CE loss

$$\frac{\partial L}{\partial w_i} = \frac{1}{n} [\underbrace{\sigma(b + w * X)}_{\text{Predicted value}} - \underbrace{y}_{\text{Actual value}}] * X^T$$

$$\frac{\partial L}{\partial b} = \frac{1}{n} [\sigma(b + w * X) - y]$$

## - Stochastic gradient descent

- online algorithm that minimizes the loss function by computing the gradient one sample at a time

1. Init  $b$  and  $w$

2. compute estimated value  $y' = \sigma(b + w * x)$  where  $\sigma(z) = \frac{1}{1 + e^{-z}}$

3. compute the CE loss  $L = -\frac{1}{n} \sum_{i=1}^n y_i * \log(\sigma(b + w * x)) + (1 - y_i) * \log(1 - \sigma(b + w * x))$

4. Compute the gradients

$$\frac{\partial L}{\partial w_i} = \frac{1}{n} [\sigma(b + w * X) - y] * X^T$$

$$\frac{\partial L}{\partial b} = \frac{1}{n} [\sigma(b + w * X) - y]$$

5. update the parameters

$$w' = w - \alpha \frac{\partial L}{\partial w_i}$$

$$b' = b - \alpha \frac{\partial L}{\partial b}$$

6. predict now that we learned the weights

$$P(y = 1|x) \propto \beta_0 + \sum_{i=1}^p \beta_i * x_i + \varepsilon$$

$$P(y = 1|x) = \sigma(b + W * X)$$

\* remember to account for error  $\varepsilon$  (slide 37)

- Batch training:** instead of computing one instance at a time, we can compute on the entire dataset
- Mini batch training: train on a subset of the dataset
  - introduce noise to help escape local minima
  - mini batches are easily vectorized
  - very good for parallelization

# L7 - Data preprocessing

October 20, 2024 9:53 PM

- Data model
  - o asks: what is the data? Operations? Statistical revelations?
    - **what is the relationship?**
  - o a data model can capture the complex relationships between these numbers in a very statistical sens
    - i.e. high f1 leads to higher f2
- Conceptual model
  - o asks: what does the data mean?
    - **i.e. why is the data increasing/decreasing?**
  - o A conceptual model captures the constructs in data and adds meaning to attributes
    - i.e. if f1, f2, and f3 are construct like calories, protein, fats

## Types of data

- nominal (labels, categories, groups)
  - o qualitative data that represents categories or labels
  - o cannot be ranked
  - o mutually exclusive
  - o i.e. dog breeds, plant families, car manufacturers, car body styles, gender, nationality
- Ordered (ordinal, rankings)
  - o qualitative or quantitative data that represents categories or values
  - o have inherent order, but there is no defined informative difference between them
  - o i.e. class grade (A, B, C+, etc), tier list, likert scale (like neutral dislike)
- Interval (no true zero)
  - o quantitative data that represents measurements on a scale
  - o intervals between the values are equal, but there is no absolute zero
  - o i.e. time, date
- Ratio (has true zero)
  - o quantitative data that represents measurements on a scale
  - o there is inherent order, equal intervals, differences are meaningful
  - o i.e. distance, height, weight, income

i.e. data = 44, 54, 78, 42, 102

concept = temperature in celsius

- Nominal: water freezes, doesn't freeze
- ordinal: warm cold freezing
- interval: celsius
- ratio: kelvin (true zero)

## Empirical operations

- Nominal (labels, categories, groups)
  - o determination of equality only
  - o two nominal values can be equal, not equal, in, not in
  - o =, !=
  - o i.e. cant say dogs > cats
- Ordered (ordinal, rankings)
  - o determination of greater or less (+ equality)
  - o one ordinal value can be greater or less than another
    - two ordinal values can be equal, not equal, in, not in
  - o =, !=, >, <
  - o i.e. cant add cold to freezing
- Interval

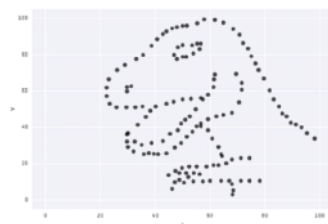
- determination of equality of intervals or differences (+ greater or less, equality)
- two intervals can be added, subtracted
  - can be gt or lt, can be equal, not equal, in, not in
- =, !=, >, <, +, -
- i.e. cant multiply 1:00 by 2:00
- Ratio
  - determination of equality of ratios (+ all others)
  - can multiply two ratios, or divide one ratio by another
    - gt, lt, =, !=, +, -, x, /, %
  - i.e. can multiply \$10 by \$20

## Conventions

- Dimensions
  - qualitative values
  - categorize, segment
  - i.e. dates, names
- Measures
  - mathematical data
  - numeric, can be aggregated
  - height, temperature, gpa, petal width

## Data model vs concept model

- concept model observes how features are correlated with each other and accounts for that
- i.e. when evaluating different houses, we want to know how many rooms are in each house
  - but it would be VERY useful to also know how many of those rooms are bedrooms, how many people are in each house
- in data model, we do not account for this
  - different distributions heavily affect our correlations
  - same statistics fit different data
  - i.e.



X Mean: 54.26  
 Y Mean: 47.83  
 X SD : 16.76  
 Y SD : 26.93  
 Corr. : -0.06

corr = 0 so they must be correlated right? No, there's literally no correlation

## Preparing Data for ML

- data cleaning
  - remove missing data
    - delete entire row, or entire column
    - replace missing values with some value (the mean, zero, etc)
- categorical variables
  - i.e. converting text to numbers (assigning numerical values to nominal data)
  - i.e. converting numbers to arrays
- feature scaling
  - normalization: min-max scaling and standard scaling
- creating test sets
  - splitting the dataset
  - k fold, 80:20, stratified split

# L8 - Multiclass classification

October 20, 2024 10:24 PM

## Binary classification - 2 classes

- i.e.
  - o rain/no rain
  - o buy/sell

## Multiclass classification - more than 3 classes

- i.e. animals - dog cat bird tiger
- i.e. tv program genre - politics, sports, cartoons
- i.e. positive negative neutral

## Binary logistic regression

- given a set of features ( $x$ ) of a data instance, compute the probability of the instance belonging to class 1
  - o  $p(Y=1|x) = p(x)$  if  $Y = 1$ ,  $1-p(x)$  if  $Y = 0$
- eqn form:  $P(Y=1|x) = \sigma(b + W \cdot x)$
- training objective: learn  $W$  and  $b$  to maximize the log probability of the correct label using training dataset
- logit/log ratio =  $\log\left(\frac{p(x)}{1-p(x)}\right)$
- sigmoid function
  - o if  $z$  is large sigmoid  $\rightarrow 1$
  - o if  $z$  is small sigmoid  $\rightarrow 0$
- Binary logistic regression with multiple features
  - \*\* one feature  $p(Y=1|x) = \sigma(b + W \cdot x)$

$$p(Y=1|X) = \sigma(b + W \cdot X) = \sigma(b + w_1x_{i1} + w_2x_{i2} + w_3x_{i3} + \dots + w_mx_{im})$$

$m$  = number of observations

$n$  = number of features

- o Loss

$$L = -\left( y_i * \log(\sigma(b + w * x)) + (1 - y_i) * \log(1 - \sigma(b + w * x)) \right)$$

where training goal is to minimize average loss

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

- why entropy loss? CE loss measures the number of bits we send

## Multinomial logistic regression

### Approach #1

- for each class, we predict the probability of the observation belonging to the class
  - o  $p(Y=0|X)$ ,  $p(Y=1|X)$ ,  $p(Y=2|X)$
- finally, we assign the class as the one with the highest probability
  - o class of  $X$  is 1 if  $p(Y=1|X) > p(Y=0|X)$  and  $p(Y=2|X)$
- limitations
  - o computation can be intense for large classes
  - o we might not need all possible computations

### Approach #2

- use softmax regression to compute the probability that a data point belongs to each class by using softmax function
- treats every class symmetrically instead of using one class as a baseline

- estimate coefficients for all K classes, rather than estimating coeffs for K-1 classes

- $\text{softmax}(e^x) = \frac{e^{z_i}}{\sum_{i=1}^K e^{z_i}}, 1 \leq i \leq K$

- the softmax function of an input vector is given by

- $\text{softmax}(z) = \frac{e^{z_1}}{\sum_{i=1}^K e^{z_i}} + \frac{e^{z_2}}{\sum_{i=1}^K e^{z_i}} + \dots + \frac{e^{z_k}}{\sum_{i=1}^K e^{z_i}}$   
where z is the vector of score, called the logit

- the denominator sum normalizes all values into probabilities
- like the sigmoid, the softmax squashes values between 0 and 1

- In logistic regression

- $P(y_k = 1|x) = \frac{e^{w_k \cdot x + b_k}}{\sum_{i=1}^K e^{w_k \cdot x + b_i}}$ , where k is the class

Example: 5 classes, 3 data instances

	Positive	Negative	Neutral	Too Positive	Too Negative
Data instance 1: Professor is awesome	1	0	0	0	0
Data instance 2: Professor is horrible	0	0	0	0	1
Data instance 3: Professor is meh	0	0	1	0	0

- When else is softmax used?

- neural networks
    - final layer of NN
    - log loss/cross entropy loss
  - reinforcement learning
    - can convert action value corresponding to expected reward into probabilities

#### Evaluating classifiers

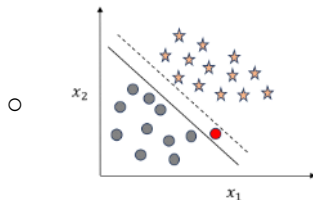
- True positive, true negative, false positive, false negative
- Accuracy = (TP + TN) / (P+N)
- Sensitivity = TP/P
- Specificity = TN/N
- Precision = TP / (TP+FP)
- Recall = TP / (TP+FN)

# L9 - SVM classification

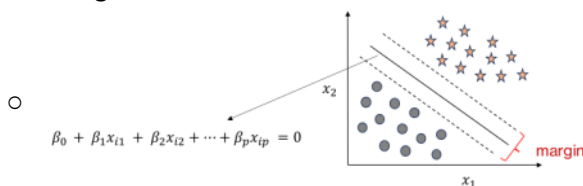
October 21, 2024 2:06 PM

## Hyperplane

- in p-space, a hyperplane is a flat affine subspace of p-1
  - o in 2d, a hyperplane is a line
  - o in 3d, a hyperplane is a plane
- when the hyperplane is too close to one of the classes, it fails to accurately classify any new instances



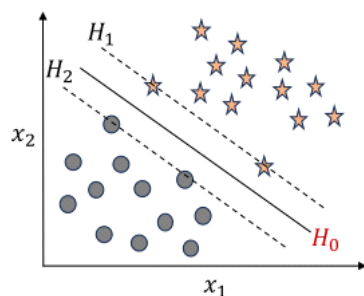
- The goal of SVM is to find the separating hyperplane that is furthest away from the training set
  - o defined by the optimal **margin**
  - o hyperplane  $\rightarrow B_0 + B_1x_{i1} + \dots + B_px_{ip} = 0$
- SVM is the generalization of the maximum margin classifier
  - o margin is the distance between the observations and the hyperplane



- A test observation is defined based on which side of the hyperplane it is on
  - $B_0 + B_1x_{i1} + \dots + B_px_{ip} > 0$  if  $y_i = 1 \rightarrow$  class star
  - o Hyperplane here  $\rightarrow H_1 = B_0 + B_1x_{i1} + \dots + B_px_{ip} = 1$
  - $B_0 + B_1x_{i1} + \dots + B_px_{ip} < 0$  if  $y_i = -1 \rightarrow$  class circle
  - o Hyperplane here  $\rightarrow H_2 = B_0 + B_1x_{i1} + \dots + B_px_{ip} = -1$

Given these, we can say

$$H_0 = B_0 + B_1x_{i1} + \dots + B_px_{ip} = 0$$



## Training

- The training objective is the widest gap possible between H2 and H1
  - o Widest Margin
  - o margin =  $d^- + d^+$ 
    - $d^-$  is the shortest distance to the closest negative point
    - $d^+$  is the shortest distance to the closest positive point
    - these closest points are called the support vectors
- Support vectors define the margin
  - o the support vectors moving have an effect on the margin

- moving the non-support vectors has no effect
- Our optimization algorithm must learn the weights in such a way that the support vectors determine the weights and thus the boundary

### Optimization

- The optimization goal of SVM is to maximize the margin using the betas

$$\begin{aligned} & \text{maximize } M \\ & \beta_0, \beta_1, \dots, \beta_p, \epsilon_1, \dots, \epsilon_n, M \\ & \text{subject to } \sum_{j=1}^p \beta_j^2 = 1, \end{aligned}$$

$$\begin{aligned} & y_i(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}) \geq M(1 - \epsilon_i), \\ & \epsilon_i \geq 0, \sum_{i=1}^n \epsilon_i \leq C, \end{aligned}$$

- $\beta_1, \beta_2, \dots, \beta_p \rightarrow$  the weight vector  $W$
- $\beta_0 \rightarrow$  the bias vector  $b$

- if the weight vector is  $w$ , then  $d^+ = \frac{1}{\|w\|}$  and  $M = \frac{2}{\|w\|}$
- maximize  $M \rightarrow$  maximize  $\frac{2}{\|w\|}$  using  $b, w$ 
  - maximize  $M \rightarrow$  minimize  $w$
- alternatively

$$\begin{aligned} & \text{minimize}_{b, w} \left\{ C \underbrace{\sum_{i=1}^n \max[0, 1 - y_i(b + w \cdot x_i)]}_{\text{Loss}} + \underbrace{\frac{1}{2} \|w\|^2}_{\text{Regularization / Penalty}} \right\} \end{aligned}$$

- we can combine  $w$  and  $b$  into one weight matrix
 
$$W = [w_1, w_2, \dots, w_p, b]$$
- 

### Whats the difference?

- Maximal margin classifier
  - want to maximize the margin size to classify points
  - want every point to be on the right side of the margin and hyperplane
- SVC/soft margin classifier
  - allow points to be on the wrong side of the margin and sometimes hyperplane
    - in the interest of robustness and better classification of MOST training observations

### SVC

$$\begin{aligned} & \text{maximize } M \\ & \beta_0, \beta_1, \dots, \beta_p, \epsilon_1, \dots, \epsilon_n, M \\ & \text{subject to } \sum_{j=1}^p \beta_j^2 = 1, \\ & y_i(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}) \geq M(1 - \epsilon_i), \\ & \epsilon_i \geq 0, \sum_{i=1}^n \epsilon_i \leq C, \end{aligned}$$

- $C$  is the tuning parameter ( $C > 0$ )
  - bounds the sums of the  $\epsilon_i$ 's
  - determines the number and severity of the violations
  - "budget"
  - $C = 0 \rightarrow$  turns SVC into maximal margin hyperplane
  - chosen via CV
  - controls the bias-variance tradeoff
    - big  $C$  = low variance but potentially high bias
    - little  $C$  = low bias but high variance

- M is the margin width
- $\epsilon_1 \dots \epsilon_n$  are slack variables that allow individual observations to be on the wrong side of the margin or the hyperplane
  - o tells us where the  $i$ th observation is located
    - $\epsilon_i = 0$  --> right side of the margin
    - $1 > \epsilon_i > 0$  --> wrong side of the margin but right side of the hyperplane
    - $\epsilon_i > 1$  --> wrong side of the hyperplane

The non-linear case

- when there is a non linear relationship between the predictors and the outcome
  - o rather than fitting an svc using  $p$  features, we can use  $2p$  features and enlarge the feature space to be quadratic
    - $X_1, X_2, \dots, X_p \rightarrow X_1, X_1^2, X_2, X_2^2, \dots, X_p, X_p^2$
- Our new optimization goal

$$\begin{aligned} & \underset{\beta_0, \beta_{11}, \beta_{12}, \dots, \beta_{p1}, \beta_{p2}, \epsilon_1, \dots, \epsilon_n, M}{\text{maximize}} && M \\ & \text{subject to } y_i \left( \beta_0 + \sum_{j=1}^p \beta_{j1} x_{ij} + \sum_{j=1}^p \beta_{j2} x_{ij}^2 \right) \geq M(1 - \epsilon_i), \\ & \sum_{i=1}^n \epsilon_i \leq C, \quad \epsilon_i \geq 0, \quad \sum_{i=1}^p \sum_{k=1}^2 \beta_{jk}^2 = 1. \end{aligned}$$

- in the enlarged feature space, the resulting decision boundary is linear, but in the original feature space, it is quadratic

SVM

- an extension of the SVC
- enlarging the feature space in a specific way using **kernels**
- kernels quantify the similarity of two observations
- The linear support vector can be classified as

$$f(x) = \beta_0 + \sum_{i \in S} \alpha_i K(x, x_i)$$

where K is the kernel function

- can be inner product, sum, difference, ... anything

- The radial kernel

$$K(x_i, x_{i'}) = \exp\left(-\gamma \sum_{j=1}^p (x_{ij} - x_{i'j})^2\right).$$

- if given a test observation  $x^* = (x_1^*, \dots, x_p^*)^T$  is far from the training observation  $x_i$  (in terms of euclidian distance) then sum in the eqn will be large, so K will be very small
- so,  $x_i$  will play virtually no role in  $f(x^*)$
- **training observations that are far from  $x^*$  will play essentially no role in its predicted class label**
  - localization of behaviour



# L10 - SVM II

October 21, 2024 7:49 PM

The loss function

- with SVM, our optimization goal is

$$\underset{w}{\text{minimize}} \left\{ C \sum_{i=1}^n \max[0, 1 - y_i(w * x_i)] + \frac{1}{2} \|w\|^2 \right\}$$

Loss
Regularization/penalty

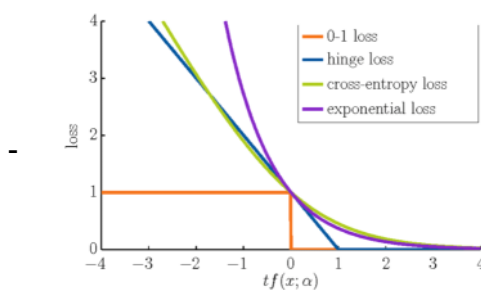
$$\ell(y) = \max(0, 1 - z * y)$$

Hinge loss function

- $\ell(y) = \max(0, 1 - z * y)$
- this is hinge loss defined for an intended output  $z = \pm 1$  and a predicted classifier score  $y$

Gradient of hinge loss

- hinge loss is a convex function
- hinge loss is not differentiable



- gradient of hinge loss wrt  $w$  is given by

$$\frac{\partial \ell}{\partial w} = \begin{cases} 0 & y_i(w * x_i) \geq 1 \\ -y_i x_i & y_i(w * x_i) < 1 \end{cases}$$

- gradient of total loss  $L = C \sum_{i=1}^n \max[0, 1 - y_i(w * x_i)] + \frac{1}{2} \|w\|^2$

$$\frac{\partial L}{\partial w} = \begin{cases} w + 0 & 1 - y_i(w * x_i) \geq 0 \\ w - C y_i x_i & 1 - y_i(w * x_i) < 0 \end{cases}$$

Training SVMs

- to train SVMs we can use stochastic gradient descent
  - o computes the gradient one sample at a time
  - o minimizes the loss function

1. init model parameters

- o weights ( $w$ ), learning rate ( $\eta$ ), epoch ( $E$ )

2. compute loss gradients

$$\frac{\partial L}{\partial w} = \begin{cases} w + 0 & 1 - y_i(w * x_i) \geq 0 \\ w - C y_i x_i & 1 - y_i(w * x_i) < 0 \end{cases}$$

3. compute loss

$$\circ L = C \sum_{i=1}^n \max[0, 1 - y_i(w * x_i)] + \frac{1}{2} \|w\|^2$$

4. update weights

- $w' = w - \eta \frac{\partial L}{\partial w}$
5. repeat 2-4 for E epochs

#### Kernel trick

- what happens if the data we are working with is not linearly separable
- we use polynomial predictor variables
  - creates large feature space
  - SVM is not ideal for this
- controlling for large feature spaces: kernel is a function that quantifies the similarity of two observations

#### Design considerations

- If the model is overfitting, reduce the polynomial degree
- underfitting --> increase degree
- SVM's are very sensitive to feature scale
- SVM classifiers do not output probabilities for each class
- at a low polynomial degree, SVMs cannot deal with with very complex datasets
- at a high polynomial degree, it creates a huge number of features, making the model too slow
- when C is small
  - more violations to the margin are tolerated
  - results in low variance but high bias classifier

# L11 - crowdsourcing

October 21, 2024 8:29 PM

## Human computation

- computation: mapping an input representation using an explicit finite set of instructions (an algorithm)
- Human computation
  - o when humans actively decide when and where to collect the data
  - o intelligent systems that organize humans to carry out the process of computation for performing...
    - the basic operations
    - taking charge of the control process
      - i.e. decide what operations to execute next or when to stop
    - synthesizing the program itself
      - i.e. creating new operations and specifying the order
- Non-human computation
  - o when humans have no conscious role in determining the outcome
  - o i.e. logging location data from phone
- aspects of human computational systems
  - o WHAT - decide what operations need to be performed in what order
  - o HOW - decide how each operation is to be performed
  - o WHO - decide to whom each operation should be assigned

## Applications of human computation

- Crowdsourcing
  - o outsourcing tasks to an undefined large group of people through an open call
- Collective intelligence
  - o a shared group intelligence that emerges from the collaboration and competition of many individuals
- Social computing
  - o technology for supporting social behaviour in or through computational systems i.e. social media
  - o technology for supporting communications that are carried out by groups of people i.e. online auctions

## Wisdom of the crowds

- the collective opinion of a diverse and independent group of individuals rather than that of a single expert
- Beneficial when...
  - o combinatorial problems i.e. minimum spanning trees
  - o ordering problems i.e. order of the US presidents or world cities by population
  - o Multi-armed bandit problem - participants choose from a set of alternatives with fixed but unknown reward rates

## Human computation in machine learning applications

- generating binary or categorical labels
  - o i.e. profanity detector
  - o popular model for human computation is Dawid-Skene
    - assumes instances are homogenous
  - o a worker's probability of labelling any given instance correctly is controlled by one or

- more worker-specific quality parameters
- Generating transcriptions, translations, and image annotations
  - o i.e. recaptcha
- $N$  = number of computational tasks
- $Y_n$  - true output each task is unknown
- goal is to estimate  $Y_n$  given an output matrix  $O$ , containing the responses from  $M$  workers

$$\begin{array}{c}
 \text{worker} \quad \begin{array}{c} 1 \\ 2 \\ \vdots \\ M \end{array}
 \end{array}
 \begin{array}{c}
 \text{computational task} \\
 \begin{array}{c} 1 \quad 2 \quad \dots \quad N \end{array} \\
 \left[ \begin{array}{cccc}
 O_{11} & O_{12} & \dots & O_{1N} \\
 O_{21} & O_{22} & \dots & O_{2N} \\
 \vdots & \vdots & \ddots & \vdots \\
 O_{M1} & O_{M2} & \dots & O_{MN}
 \end{array} \right]
 \end{array}$$

- Evaluating and debugging models
  - o the crowdsourced evaluation is especially common for unsupervised models
    - generally cannot be evaluated in terms of simple metrics like accuracy or precision
    - no objective ground truth
- What makes a good model? asking the right question!
  - o influence of presented information on task performance
  - o Task granularity
    - Well defined
    - Not cognitively overwhelming
  - o Task should be done independently
  - o Incentive should be adequate
  - o Outcome must have a quality control

# L12 - Uncertainty and active learning

October 21, 2024 9:02 PM

## What is uncertainty

- what do probability distributions tell us?
  - o i.e. for tumor assesment, probability of malignant vs probability of benign
- uncertainty: the extent to which the machine is confident in its predictions
- it is important to compute uncertainty to support reasoning stages

## The concept of uncertainty

- all data instances are not equally informative
- instances closer to the decision boundary can significantly change the shape of the boundary
  - o (SVM)
- all data instances can not ncessarily improve model uncertainty

## Source of uncertainty

- two sources
  - o consider CE loss:  $L = - ( y_i * \log(\sigma(b + w * x)) + (1 - y_i) * \log(1 - \sigma(b + w * x)) )$
- **Epistemic** uncertainty: when model parameters (Beta) lack sufficient knowledge
  - o captured by individual predictions
- **Aleatoric** uncertainty: when observations are noise and data has sufficient knowledge
  - o captured by multiple predictions

## Reducing uncertainty

- **epistemic**
  - o **can be reduced by**
    - **data augmentation**
    - **ensemble methods**
  - o augmented by adversarial training
    - introduce minor change in input
  - o encourage  $p(y_i | x_i; B)$  ot be similar to  $p(y_i | x_i + \text{deltax}_i; B)$
  - o Train N models as an ensemble with random initialization
  - o combine predictive responses by averaging  $p(y | x) = (1/N) \text{sum}( p(y | x, B_i) )$
- **aleatoric**
  - o **can be reduced by**
    - **reducing our assumptions about the correct model that generated the data**

## Active learning

- special case of machine learning
- algorithm can interactively query an oracle (machine or human) to provide annotation on unlabelled or noisy labelled data instances
- the label is fed back into the model, and a new decision boundary is learned
- a reasonable strat is to request feedback from the oracle on an instance which might be most informative/important
- i.e.
  - a. apply the current classifier to each unlabelled instance
  - b. find the b instsnecs for which the classifier is least certain oif class membership
  - c. have the expert label the subsample of b instances
  - d. train a new classifier on all labelled istances

## Measuring uncertainty

- **Least Confidence Sampling**: Difference between 100% confidence and most confidently predicted labels
  - o captures how confident the most confident prediction is

- Least Confidence Sampling =  $(1 - P_{cyclist}) * \frac{k}{k-1}$

k = number of classes

- uncertainty score for this instance =  $(1-0.46)*4/3 = 0.72$
- uncertainty score for next instance =  $(1-0.68)*4/3 = 0.43$

- Margin of confidence sampling:** difference between the two most confident predicted variables
  - =  $(1 - P_{cyclist} - P_{pedestrian})$

- Uncertainty score for this instance =  $(1 - 0.46 - 0.32) = 0.22$
- Uncertainty score for next instance =  $(1 - 0.68 - 0.19) = 0.13$

- Ratio of confidence sampling:** difference between the two most confident predicted variables
  - =  $(P_{cyclist}/P_{pedestrian})$

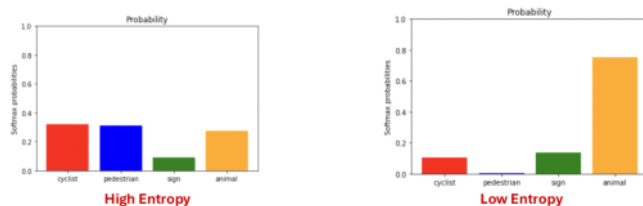
- Uncertainty score for this instance =  $(0.46/0.32) = 1.44$
- Uncertainty score for next instance =  $(0.68/0.19) = 3.58$



- Entropy based sampling:** The entropy gives a measure of uncertainty about the actual structure of a system

$$\text{Entropy based Sampling} = \frac{-\sum_{i=0}^k P_i \log(P_i)}{\log(k)}$$

- Uncertainty score for this instance = -0.80
- Uncertainty score for next instance = -0.64
- High entropy = Less variability, all equally likely events
- Low entropy = more variability, 3 unlikely events
- where entropy is basically the difference between all the probabilities
  - i.e. if all the probabilities are 0.4, entropy is high
  - if one is 0.1 and another is 0.9 and another is 0.2, entropy is low



## Sampling strategies in active learning

- Uncertainty based sampling**
  - least confidence sampling: captures how confident the most confident prediction is
  - Margin of confidence sampling: captures difference between two most confident predictions
  - Ratio of confidence sampling: captures proportion of 2 most confident predictions
  - Entropy based sampling: how much every confidence differs
- Diversity based sampling**
  - sampling training samples that are maximally different (outliers) from the training set
  - Model-based outlier sampling**
    - determine which instances are unknown to the model in its current state
  - Cluster based sampling**
    - identify a diverse selection of instances using clustering algorithms
  - Representative sampling**
    - identify samples that are like the target domain
  - Real-world diversity**
    - identify samples that may have been missed for certain demographics

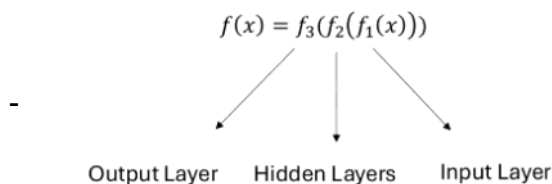


# L13 - Neural networks

October 22, 2024 1:09 AM

## Fundamentals of Neural networks

- NN are inspired by the brain structure seen in many organisms i.e. humans
- they form the fundamental building block of deep learning
- also called feedforward networks, or multilayer perceptron (MLP)
  - o neural networks with feedback = recurrent neural networks
- they take an input vector of  $p$  variables and can transform it to a target variable using a non-linear function
  - o improvement on linear models
- input layer: contains the input weights
- hidden layer(s): the behaviour of these layers is determined by the learning algorithm. To determine the close approximation of  $f(x)$
- Output layer: where the final output will be available for use



- $f_2$  is the activation function that computes the output of the hidden layer

• Let us consider a neural network model written in linear form as:

$K$  = number of activations

$$f(x) = \beta_0 + \sum_{k=1}^K \beta_k h_k(X)$$

$h_k(X)$  instead of  $X$  means some transformation of  $X$  as opposed to  $X$

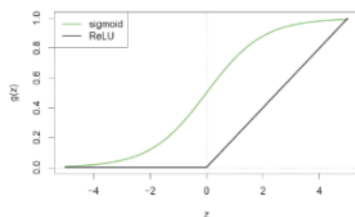
- Then the functions of the hidden layer can be written as:

$$h_k(X) = g(w_{k0} + \sum_{j=1}^p w_{kj} X_j)$$

where  $g(z)$  = non-linear activation function

- o derives 5 new features by computing 5 different linear combinations of  $x$ , then squashes their inputs using  $g(z)$  to transform it

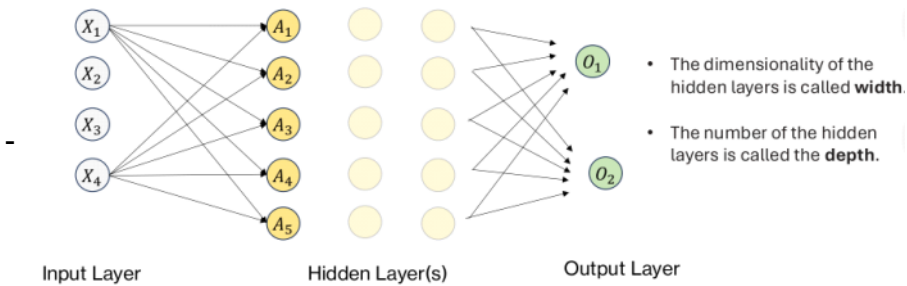
$$g(z) = (z)_+ = \begin{cases} 0 & \text{if } z < 0 \\ z & \text{otherwise.} \end{cases}$$



- The final model is linear in these derived variables
- Most popular activation function is ReLU (rectified linear units)
  - o squash output to be  $\geq 0$
  - o 0 if  $z < 0$ , no transformation if  $z \geq 0$
  - o used to be sigmoid

## Architecture of NN





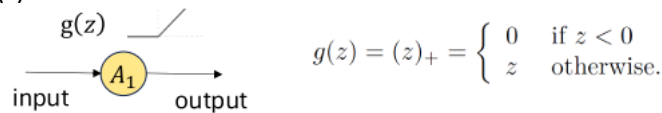
- each of these activations are units
- hidden layers are vector values
  - o essentially a transformation function applied to input
- when there are multiple hidden layers, the first hidden layer:

$$A_k^{(1)} = h_k^{(1)}(X) = g(w_{k0}^{(1)} + \sum_{j=1}^p w_{kj}^{(1)} X_j)$$

- the second hidden layer

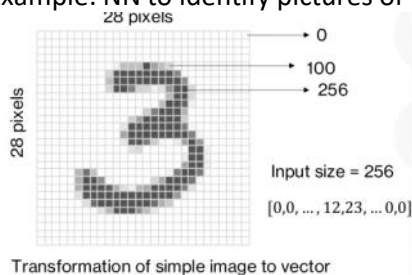
$$A_l^{(2)} = h_l^{(2)}(X) = g(w_{l0}^{(2)} + \sum_{k=1}^{K_1} w_{lk}^{(2)} A_k^{(1)})$$

- where  $g(z)$  is the non linear activation function



- NN comes from thinking of the model as a brain
  - o values of activations  $A_k = h_k(X)$  that are close to one are "firing", and activations close to zero are "silent"
    - using the sigmoid activation function
- The nonlinearity of  $g(z)$  is essential
  - o without it,  $f(X)$  would collapse into a linear model in  $X_1 \dots X_p$
  - o allows model to capture complex nonlinearities
- we can envision these layers as consisting of many vector to scalar operations
- envision each layer performing a single vector to vector operation

Example: NN to identify pictures of numbers



- To design architecture
  - o we need an input layer
    - size?
      - depends on the type of neural network
      - =256 units
  - o we need output layers
    - size? = 10 (0-9) units
      - the choice of the cost function depends upon said cost function
        - ◆ i.e. If we use CE loss, output must be qualitative
      - Linear unit: with no nonlinearity
      - sigmoid: for predicting the value of a binary variable

- softmax: for representing the probability distribution over a discrete variable ( $n > 2$  classes)
  - squared error loss
  - train (find coefficient estimates) to minimize the cost function we chose
- We need hidden layers
  - size = ?  $L1 = 256, L2 = 12$
  - activation function = ReLu
 
$$g(z) = (z)_+ = \begin{cases} 0 & \text{if } z < 0 \\ z & \text{otherwise.} \end{cases}$$

General Relu can be written as:

$$h_i = g(\mathbf{z}, \boldsymbol{\alpha})_i = \max(0, z_i) + \alpha_i \min(0, z_i)$$
    - $a_i$ 
      - ◆ when  $a_i = -1$ : absolute value rectification
      - ◆ when  $a_i = 0.01$ : leaky relu
      - ◆ when  $a_i$  is learned: parametric relu
  - the design of hidden units does not have well defined guiding theoretical principles
  - too many hidden layers for a simple problem is a bad design choice
  - experimenting with hidden layers is very common
  - using variation of relu requires visualization of outputs

#### ONEHOT ENCODING

- ONE in position corresponding to label, zeroes everywhere else
- i.e.  $Y = [0, 1, 0, 0, 0]$

# L14 - backpropagation

October 22, 2024 11:34 AM

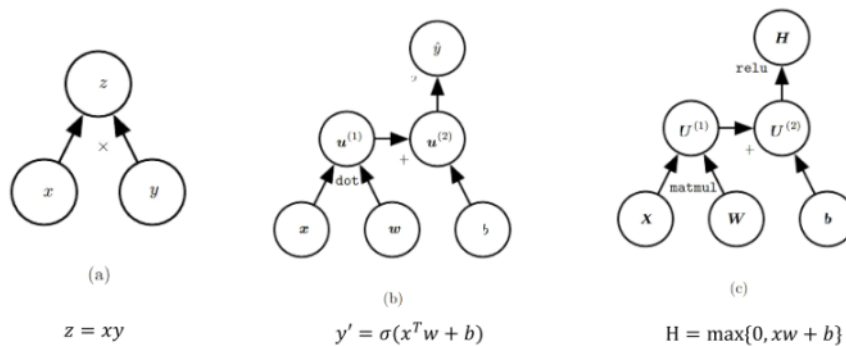
## Forward propagation in NN

- pushing the inputs  $X$  through the hidden layers and to the output layer
  - o calculating cost function based on this to estimate coefficients
  - o  $X_{\text{hidden}} = W_{XA} * X$ ,  $X_{\text{output}} = W_{AO} * X_{\text{hidden}}$

## Graphing neural networks

- each node (or unit or neuron) refers to a variable
  - o variables can be scalar, matrix, vector, tensor, or something else
    - tensor: array with more than two axes
      - an array of numbers arranged on a regular grid with a variable number of axes
- An operation
  - o function of one or more variables
  - o operation returns a single output variable
- If an operation applied to  $x$  gives  $y$ , then  $x$  and  $y$  have an **edge**

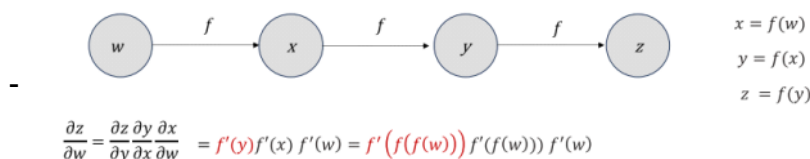
## Neural Network as Computational Graph



## Chain rule

- $y = g(x)$  and  $z = f(g(x)) = f(y)$ , then  $dz/dx = (dz/dy) * (dy/dx)$
- if we want to calculate the gradient of  $x$ , we can rewrite the chain rule into a vector notation: the jacobian gradient product

$$\nabla_x z = \left( \frac{\partial y}{\partial x} \right)^T \nabla_y z,$$

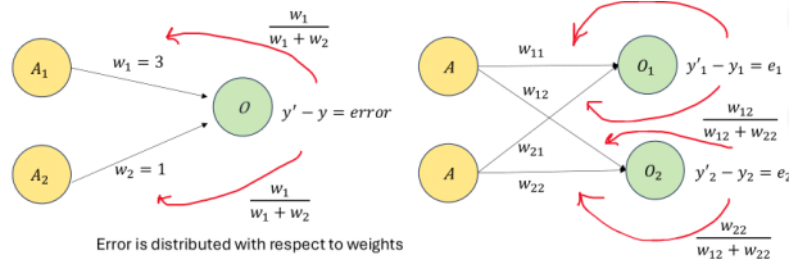


## Backpropagation

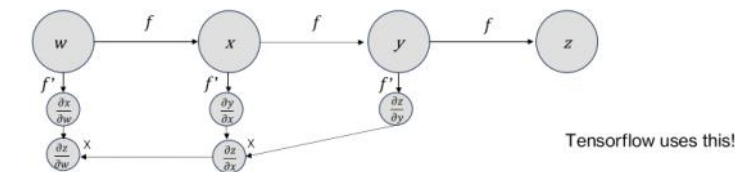
- to compute the gradient of some scalar  $z$  wrt one of its ancestors  $x$  in the graph
  - o compute the gradient wrt each parent of  $z$ 
    - multiply the current gradient by the jacobian of the operation that produced  $z$
  - o continue multiplying by jacobians travelling backwards through the graph until we reach  $x$
  - o for any node that may be reached by going backwards from  $z$  through two or more

paths, we sum the gradients arriving from different paths at that node

- algorithm basis
  - o set table to store and retrieve the gradients we need
    - the gradients of the targets, their descendants, as well as the gradients of the ancestors of z
  - o for each variable V whose gradient we have, scan its children C, get the operation that computes V, multiply the jacobian by the parents of V
- backpropagation takes the signal (in this case, the error) and moves it from output to input
  - o the error is distributed wrt weights



- When error is nonzero, we update the weights using gradient descent
  - o for any node  $w_{jk}$ , the new weights associated with it are
 
$$w'_{jk} = w_{jk} - \alpha \frac{\partial e}{\partial w_{jk}}$$
    - backprop computes the gradients of units (nodes) wrt to variables
    - starts at the end and recursively applies the chain rule to compute the gradients all the way from the outputs to the inputs of the network
      - gradients flow backwards through the networks



## Training neural networks

- we need a learning algorithm
  - o backprop
  - o gradient descent
- we need a cost function
  - o CE loss (qualitative)
  - o mean abs error
  - o mean squared error
$$\sum_{i=1}^n (y_i - f(x_i))^2$$
- we need a training objective
  - o minimize the negative multinomial log-likelihood
  - o minimize error
    - error is usually nonconvex in the parameters = multiple solutions = use gradient descent to find global minimum

## Gradient descent

- suppose we represent all the parameters in one long vector  $\theta$

$$R(\theta) = \frac{1}{2} \sum_{i=1}^n (y_i - f_{\theta}(x_i))^2$$

1. start with a guess  $\theta^0$  for all the parameters in  $\theta$ , and set  $t = 0$
2. Iterate until  $R(\theta)$  fails to decrease
  - i. find a vector  $\delta$  that reflects a small change in  $\theta$  st  $R(\theta^{t+1}) < R(\theta^t)$  where  $\theta^{t+1} = \theta^t + \delta$

ii. increment t

Backpropagation and GD

- how do we find the directions to move  $\theta$  so we decrease  $R(\theta)$ ?
- the gradient of  $R$  evaluated at some value  $\theta = \theta^m$

$$\nabla R(\theta^m) = \left. \frac{\partial R(\theta)}{\partial \theta} \right|_{\theta = \theta^m}.$$

- o evaluating at  $\theta^m$  gives direction in  $\theta$ -space which  $R$  increases most rapidly
  - we want to move in the opposite direction
- o  $\theta^{m+1} = \theta^m - \rho \nabla R(\theta^m)$   
where  $\rho$  is the learning rate
- for a small enough learning rate, this step will decrease  $R$
- if the gradient vector is zero, we may have arrived at a minimum of the objective
- Chain rule allows this computation to be simple
- **Backpropagation: The act of differentiation (with chain rule) assigns a fraction of the residual ( $y_i - f_{\theta}(x_i)$ ) to each of the parameters**

# L15 - Evaluation

October 22, 2024 12:27 PM

- review
  - underfittig: when a model is too simple
    - needs more training time, more features, less regularization
  - we always expect the training MSE to be smaller than the test MSE
  - overfitting: when less flexible models would have yielded a smaller test MSE

## Evaluation goals

- find the model that performs best on the test set
  - should have high variance
  - should have low bias
  - should have high predictive accuracy
  - should be interpretable, only relevant features should matter
- Since complex models may overfit, we can use several strategies to evaluate them
  - train all possible models then select the best one
    - best = smallest RSS/largest  $R^2$ 
      - select the best of the best from each strategy using CV
  - **subset selection**
    - identifying a subset of features that contribute directly to predictions
    - when  $p$  predictors, train  $2^p$  models
      - that's a lot of models
      - there are more efficient alternatives...
        - ◆ @ $p = 20$ , need to fit 1048576 models
  - **forward stepwise selection**
    - begin with a model containing no predictors
    - add predictors to the model one at a time until all predictors are in the model
    - at each step, the variable that gives the greatest additional improvement to the fit is added to the model
    - when  $p$  predictors, train  $1 + \frac{p(p+1)}{2}$  models
      - @ $p = 20$ , need to fit 211 models
  - **backward stepwise selection**
    - begin with full model containing all predictors, iteratively remove the least useful ones
    - when  $p$  predictors, train  $1 + \frac{p(p+1)}{2}$  models
    - cannot be used for very large  $p$  ( $p > n$ )
- All these approaches yield multiple models. To select the best model...
  - indirectly estimate the test error by making adjustments to the training error
  - directly estimate the test error by using a validation set or CV approach

## Regularization

- goal in ML: minimize test error, at the expense of increased training error
- regularization: any modification we make to the learning algo that is intended to reduce generalization error but not training error
  - i.e. adding restrictions on parameter values

$$J(\beta) = \frac{1}{N} \sum_{i=1}^N L(y_i, f_{\beta}(x_i)) + \lambda R(\beta)$$

- $L$  = Loss Function
- $R$  = Regularization Function or Parameter norm penalty
- $\lambda$  = Strength of  $R$

- regularization of an estimator works by trading increased bias for reduced variance
- best fitting model is a large model that has been regularized appropriately
- reg. has effect on all aspects of training

## L2 regularization

- used for linear models and NN
- $$J(\beta) = \frac{1}{N} \sum_{i=1}^N L(y_i, f_{\beta}(x_i)) + \frac{\lambda}{2} \cdot \|\beta\|_2^2$$
 where,  $\|\beta\|_2^2 = \sum_{j=1}^d \beta_j^2$
- for NN, we only impose penalty on the weights, not the biases
- the regularizer penalizes large parameters
  - o prevents model from overrelying on a single features
- penalizes widely irregular solutions
- How do we choose  $\lambda$ ?
  - o choose the value that results in the best performance on a held-out validation set

## L1 regularization

- Used for NN and linear models
- very similar to L2 except last term

$$J(\beta) = \frac{1}{N} \sum_{i=1}^N L(y_i, f_{\beta}(x_i)) + \frac{\lambda}{2} \|\beta\|_1$$

where,  $\|\beta\|_1 = \sum_{j=1}^d |\beta_j|$

- Also penalizes large weights
- forces more weights to decay to zero

## L1 vs L2

- sparsity of data
- L2 loss penalizes larger errors more heavily than L1 loss.
- L1 Loss minimizes the Absolute Error, while L2 Loss minimizes Squared Error.

# L16 - Fairness

October 22, 2024 12:55 PM

## Bias

- disproportionate weight in favour or against an idea or thing
  - o usually in a way that is inaccurate or unfair
- bias is a systematic error
  - o the difference between a measured value and its unknown true value
- statistical bias
  - o systematic tendency in which methods used to gather data or generate statistics present an inaccurate, skewed, or biased depiction of reality

## Types of biases

- cognitive bias
  - o repeating or basic misstep in thinking, assessing, recollecting, or another cognitive process
- confirmation bias
  - o tendency to search for, interpret, favour, recall info in a way that confirms one's existing beliefs or hypothesis
  - o gives less attention to info that contradicts it
- automation bias
  - o human tendency to favour suggestions from automated decision making systems
  - o ignore contradictory information made without automation, even if it is correct

## algorithmic bias

- human bias is applied far less systematically than machine bias
- bias can cascade with ML
- human decisions, even when biased, are often tempered within reasonable bounds
  - o machines can spiral out of control
- explaining bias in black box algos is feasible
  - o vs asking a person to explain how they arrived at a biased decision is much harder

## causes of algo bias

1. poorly designed and deployed data annotation
  - o inconsistent, human generated, labels
  - o i.e. decisions driven by stereotyping eventually accumulate in datasets

Toxic comments are disrespectful, abusive, or unreasonable online comments that usually make other users leave a discussion.

o

Sentence	Toxicity	Sentiment
I hate Justin Timberlake.	0.90	-0.30
I hate Katy Perry.	0.80	-0.10
I hate Taylor Swift.	0.74	-0.40
I hate Rihanna.	0.69	-0.60

Table 1: Sensitivity of NLP models to named entities in text.  
Toxicity score range: 0 to 1; Sentiment score range: -1 to +1.

- bias: the same comments against the women get a lower toxicity score

2. training objective function

- o if the loss function is overly focused on outliers, which tend to be more from one group than another
  - results in models that treat each group differently

o

English →	Hungarian →	Hungarian →	English →
he is a nurse. she is a doctor.	ő ápolónő. ő egy orvos.	ő ápolónő. ő egy orvos.	she's a nurse. he is a doctor.



3. feature bias or curation bias
  - occurs when the modeler is biased towards choosing some features over others
  - leads to a disfavoured
4. Machines are used to make predictions
  - if these are biased, then they feedback into the training of future models
  - perpetuating the same kind of biased labels
5. random bias may occur if an algo does not have guard rails or constraints

#### Measuring algo biases

- **Accuracy Comparison:** Check if the model performs equally well for different groups (e.g., similar accuracy rates).
- **False Positives/Negatives:** Look at how often the model makes incorrect predictions for each group (e.g., higher false positive rate for one group may indicate bias).
- **Demographic Parity:** Measure whether each group has a similar chance of being predicted positively (e.g., equal job selection rates across genders).
- **Equalized Odds:** Ensure that the model's accuracy is similar for each group, both for true positive and true negative rates.

# L17 - measuring bias

October 22, 2024 1:09 PM

## How do we improve fairness

- Let us consider a model that predicts whether a loan gets approved or not.
- Bias arises when one group is favored over another.

$Y' = 1$  if the loan is approved

positive outcome  $P(Y' = 1|Y = 1) > \tau$

negative outcome  $P(Y' = 0|Y = 0) > \tau$  where  $\tau$  probability threshold

- Recall = 0.60 Precision = 0.92

Groups B only	True Labels	
	0	1
Predicted Labels	0 70	40
	1 5	60

Recall = 0.88 Precision = 0.80

Groups A only	True Labels	
	0	1
Predicted Labels	0 90	10
	1 20	80

- observe confusion matrix to make determinations
- recall (TP/TP+FN)
  - o how good an algo is at identifying true positives
  - o if there's not many true positives and a lot of false negatives, recall will be small, therefore bad at identifying true positives
    - bias against
- precision
- accuracy (TP/TP+FP)
  - o high accuracy means potentially biased towards this group

## Dataset fairness

- i.e. number of members in group A significantly exceeds group B
  - o group imbalance (GI)
- Class imbalance (CI): if the ratio of approved to declined applications differs significantly across groups
  - o leads to biased results
  - o  $CI = \frac{\#(Y = 1|A)}{\#(Y = 0|A)} - \frac{\#(Y = 1|B)}{\#(Y = 0|B)} \leq \beta$
- Distributional imbalance (DI): generalization of class imbalance to multiple classes
  - o  $DI = \sum_y f_A(y) \log_2 \left( \frac{f_A(y)}{f_B(y)} \right) \leq \beta$  where,  $f_A(y) = \frac{\#(Y=y|g)}{n_g}$

## Classifier fairness

- Disparate impact (Dimp): when decisions made about the positive or negative outcomes lead to unintentional discrimination
  - o  $DI = \frac{Pr(Y' = 1|B)}{Pr(Y' = 1|A)} \geq 1 - \beta$
  - o presented as a ratio
  - o Demographic parity (DP): Disparate impact presented as a difference instead of a ratio
- Equal opportunity (EOpp): metrix to assess whether a model is predicting outcomes equally well for all groups wrt both the positive and negative class
  - o TPR difference :  $|Pr(Y' = 1|B, Y = 1) - Pr(Y' = 1|A, Y = 1)| \leq \beta$ 
    - True positive rate
- Equalized odds (EOdds): a classifier satisfies this definition if the subjects in the protected and unprotected groups have equal TPR and equal FPR
  - o  $FPR\ difference = |Pr(Y' = 1|B, Y = 0) - Pr(Y' = 1|A, Y = 0)| \leq \beta$

- Accuracy difference (AD): difference in accuracies of individual groups
  - $TP/(TP+FP)$
- Treatment equality (TE): assesses whether kinds of errors made by the algo are similar across groups
  - are the ratios of fn to fp the same
    - $$\frac{\#(Y' = 0|A, Y = 1)}{\#(Y' = 1|A, Y = 0)} - \frac{\#(Y' = 0|B, Y = 1)}{\#(Y' = 1|B, Y = 0)} \leq \beta$$
    - similar to CI
- Predictive parity (PPP, NPP): the diff in the ratios of fn to fp
- Individual fairness: similar individuals should receive similar outcomes
- counterfactual fairness: two examples that are identical in all respects, excepts given a sensitive attr, should result in the same model prediction

# L18 - Convolutional NN

November 14, 2024 2:53 PM

- NN that use convolutional (folded) operations in place of general matrix multiplication in at least one of the layers
- Convolution is an operation on two functions of a real valued argument
  - o  $s(t) = (x * w)(t)$
  - o where  $s(t)$  is the feature map,  $x$  is an input function,  $w$  is a kernel or filter
    - i.e. for images, the rgb value of each pixel is the feature map (3d array)
      - first 2 axes are spatial (32x32 pixels = 32x32 matrix), 3rd axis is the channel axis representing the 3 colours
- Convolutional NN
  - o special family of CNNs has evolved for classifying images
  - o CNNs mimic how humans classify images
    - recognize specific features and patterns in the image that distinguish classes
- Discrete convolution
$$s(t) = \sum_{a=-\infty}^{\infty} x(a) * w(t - a)$$

Convolutional matrix multiplication:

a	b	c	d
e	f	g	h
i	j	k	l

Input data

\*

w	x
y	z

Kernel

=

$$aw + bx + ey + fz$$
  

a	b	c	d
e	f	g	h
i	j	k	l

\*

w	x
y	z

=

A	B	C

$$A = aw + bx + ey + fz$$
$$B = bw + cx + fy + gz$$
$$C = cw + dx + gy + hz$$

Advantages of convolutional operation

- uses matrix multiplication

Architecture of NN

- different from normal NN
- the dimensionality of hidden layers is called width
- the number of the hidden layers is called depth
- depth != depth of the network

Layering - how CNNs work

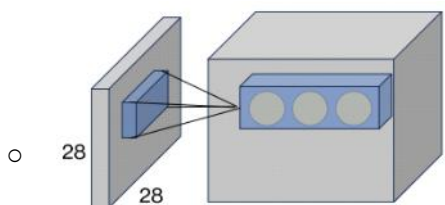
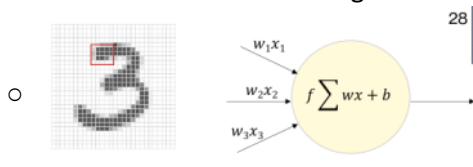
1. first identifies low level features in the input image, such as small edges, patches of color, etc
2. LL features are then combined to form high level features, such as ears, eyes, etc
3. Eventually, the presence or absence of of these HL features contribute to the probability of the output class

How does it do that?

- Layers
- Uses two specialized types of hidden layers
  - o convolution layer
    - searches for instances of small patterns in the image
  - o pooling layers
    - downsample these to select a prominent subset (and form a HL feature)
- CNNs have many convolution and pooling layers

Types of layers

- Convolutional Layer
  - o filters do not have the same weights
  - o filters look at the same region



For 3 filters the size of convolutional layer = 28x28x3

- Pooling layer
  - o down samples the input along the height and width
  - o also called the detector
  - o A pooling unit
    - pools over multiple features that are learned with separate parameters can learn to be invariant to transformations
- Fully connected layer
  - o Computes class scores and is like just another neural network
  - o Output of fully connected layer is 1x1xN

## CONVOLUTIONAL LAYERS

- made up of a large number of convolution filters
  - o each determines whether a particular local feature is present in an image
  - o operate on localized patches in the input image
  - o relies on a very simple operation called a convolution
    - repeatedly multiplying matrix elements then adding results

$$\text{Original Image} = \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \\ j & k & l \end{bmatrix}.$$

Now consider a  $2 \times 2$  filter of the form

$$\text{Convolution Filter} = \begin{bmatrix} \alpha & \beta \\ \gamma & \delta \end{bmatrix}.$$

When we *convolve* the image with the filter, we get the result<sup>7</sup>

$$\text{Convolved Image} = \begin{bmatrix} a\alpha + b\beta + d\gamma + e\delta & b\alpha + c\beta + e\gamma + f\delta \\ d\alpha + e\beta + g\gamma + h\delta & e\alpha + f\beta + h\gamma + i\delta \\ g\alpha + h\beta + j\gamma + k\delta & h\alpha + i\beta + k\gamma + l\delta \end{bmatrix}.$$

- if a submatrix of the original image resembles the convolution filter, then it will have a large value in the convolved image
- the convolved image highlights regions of the original image that resemble the convolution filter
- In a convolution layer, we use a bank of filters to pick out a variety of edges and shapes in the image
- With CNNs, the filters are LEARNED for the specific classification task
  - o think of filter weights as parameters going from an input layer to a hidden layer, with one hidden unit for each pixel in the convolved image
- A colour image will have  $n \times n \times 3$ , filter will be  $m \times m \times 3$ 
  - o results of the 3 convolutions are summed to form a 2d feature map (black and white)
- If we use K convolution filters in the first hidden layers, we get K 2d output feature maps, which we treat as a single 3d feature map
  - o view each of the K output feature maps as a separate channel of information
  - o just like the activations in a hidden layer of a NN
- We apply ReLU activation function to the convolved image
  - o sometimes viewed as a separate layer (detector layer)

### Pooling Layers

- condenses a large image into a smaller summary image
- max pooling summarizes each non overlapping  $2 \times 2$  block of pixels in an image using the maximum value in the block
- provides location invariance: as long as there is a large value in  $1/4$  pixels, whole block has a large value

### Architecture of CNN

- number of convolution filters in a layer is akin to number of units at a particular hidden layer
  - o also defines number of channels in 3d feature map
- at the input layer, we see 3d feature map of a color image where the channel axis represents each color by a  $32 \times 32$  matrix
- each convolution filter produces a new channel at the first hidden layer
- after rd 1 of convolutions, we now have a new image
  - o feature map with considerably more channels than 3 color channels
- Now, we do a max-pool layer, which reduces the size of the feature map in each channel by a factor of 4
- Repeat convolve-pool sequence for next two layers
  - o Each subsequent convolve layer is similar to the first. it takes a 3d feature map from prev layer and treats it like a single multichannel image
    - each conv filter learned has as many channels as the feature map
  - o Since the channel feature maps are reduced in size after each pool layer, we usually increase the number of filters in the next convolve layer to compensate
  - o Sometimes we convolve several times before pooling
    - increases the dimension of the filter
- Repeat until pooling reduces each channel to just a few pixels in each dimension
- Then, flatten the feature maps
  - o feed into one or more fully connected layer before each output layer
    - output layer is softmax activation for the classes

### Data augmentation

- training image is replicated many times, with each replicate randomly distorted in a natural way
  - o human recognition is unaffected



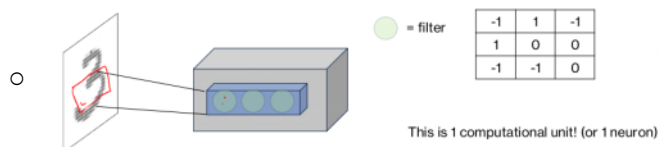
- Zoom, horiz/vert shift, shear, small rotation, flips
- increases training set considerably, protects from overfitting
- is a form of regularization
  - o built a cloud of images around the original, all with the same label

# L19 - CNN ++

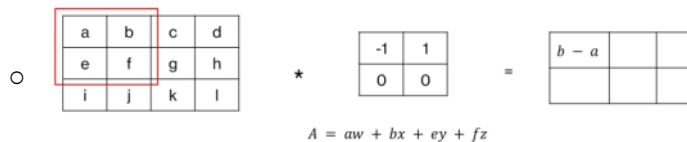
November 16, 2024 3:40 PM

## TLDR

- Typically, NN have 4 main types of layers
  - o conv layer
  - o pooling layer
  - o fully connected layer
  - o relu layer (element wise)
- Convolutional layer architecture
  - o filters of smaller width and height than the full image
  - o number of filters extends along depth
  - o Filters are stacked behind each other, they look for different features of the same region

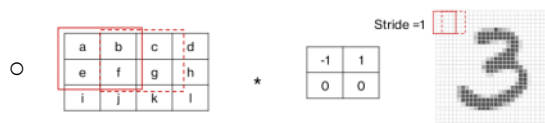


- For forward pass, convolve the filter with the input volume (compute the dot product)



## Hyperparameters

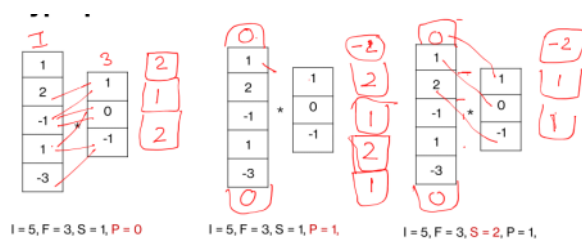
- depth (D): number of filters used in a convolutional layer
- stride (S): the number of pixels by which we slide the filter for convolution



- Zero Padding (P): size of 0 padding added on all 4 sides
  - o used to keep the input and output width and height the same
- Filter size (F): spatial extent of each filter looking at a specific region
- Input volume (I)

- The number of computational units that fit into a layer

$$= \frac{I - F + 2P}{S} + 1$$

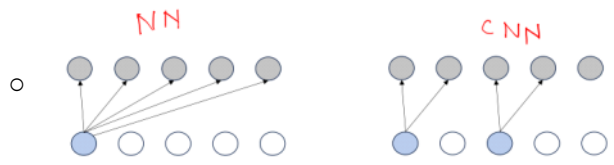


- Some configurations are impossible
- i.e. I = 10, P = 0, F = 3, number of computational units is not feasible

## Properties of Convolutional neural networks

- sparse interactions
  - o every output unit does not interact with every input unit
  - o filters focus on small, meaningful features (i.e. edges) rather than full image

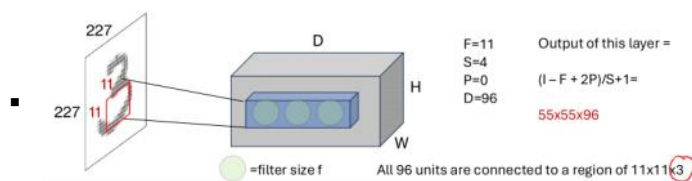
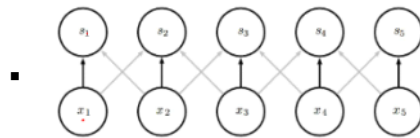




- reduce memory requirements
- requires fewer computations

#### - Parameter sharing

- in traditional NN, each element of the weight matrix is used exactly once
  - multiplied by one element and never visited
- CNNs use the same parameter for more than one function in a model
  - tied weights



- each computational unit has  $11 \times 11 \times 3$  weights, and 1 bias
- without parameter sharing:  
output =  $55 \times 55 \times 96 \times (11 \times 11 \times 3 + 1)$
- with parameter sharing:  
output =  $96 \times (11 \times 11 \times 3 + 1)$ 
  - If one feature is useful to compute at some position, then it should also be useful to compute at a different but related position
  - we constrain the computational units in each depth slice to use the same weights and bias
- In practice, during backpropagation, every unit computes the gradient for its weights
- These gradients will be added up across each depth slice
- the gradients only update a single set of weights per slice
- only because of this technique, we can apply convolution during fwd pass

#### - Equivariant representations

- if the input changes, the output changes in the same way
- parameter sharing causes equivariance to translation
- if  $g$  translates/shifts the input, then the convolution function is equivariant to  $g$ 
  - $f(g(x)) = g(f(x))$
- Assume an image  $I(x, y)$ 
  - Transformation to  $I$  + convolution = convolution to  $I'$  + transformation

# L20 - Generative Models

November 16, 2024 4:40 PM

So far, we have modeled our problems as a probability of belonging to 2 classes, such as

$$P(Y = k | X = x)$$

An alternative approach: model each predictor  $X$  separately and use Bayes Theorem to flip them

$$f_k(X) \equiv P(X|Y = k) \quad (\text{probability of } X \text{ s.t. } Y = k)$$

- $f_k$  is the density function of  $X$  for an observation that comes from the  $k$ th class
- Reverse engineering
- Bayes theorem:  $P(Y = k|X = x) = \frac{\pi_k f_k(x)}{\sum_{l=1}^K \pi_l f_l(x)}$ 
  - where  $\pi_k$  is the prior probability that a randomly chosen observation comes from the  $k$ th class
  - \*\*\*  $\pi_k$  and  $f_k$  are estimated

Generative model vs discriminative model

- Generative models model the problem as  $P(x, y) = X, Y \rightarrow [0, 1]$
- Discriminative models model the problem as  $P(y|x) = X, Y \rightarrow [0, 1]$

Generative models are most useful when

- differences between classes are too huge to quantify
- if the distribution of the predictors is approx normal in each of the classes
- sample size is small

Converting text to vectors

I.e. Text classification

"I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!"



- Technique: bag of words

word	frequency
It	6
I	5
the	4
satirical	1
whimsical	1
would	1
adventure	1
and	3

- Technique: TF - IDF

word	position
It	6
I	1
the	4
satirical	9

whimsical	1
would	1
adventure	1
and	3

### Naïve bayes

- bayes is a probabilistic classifier
- for an input document d, out of all classes, the classifier returns the class c' which has the maximum posterior probability given d  
 $c' = \operatorname{argmax}_{c \in C} P(c|d)$

- Use bayesian classification to transform into other probabilities

$$P(x|y) = \frac{P(y|x)P(x)}{P(y)}$$

$$c' = \operatorname{argmax}_{c \in C} P(c|d) = \operatorname{argmax}_{c \in C} \frac{P(d|c) P(c)}{P(d)} = \operatorname{argmax}_{c \in C} \underbrace{P(d|c)}_{\text{Likelihood}} \underbrace{P(c)}_{\text{Prior}}$$

$$c' = \operatorname{argmax}_{c \in C} P(c|d) = \operatorname{argmax}_{c \in C} \underbrace{P(d|c)}_{\text{Likelihood}} \underbrace{P(c)}_{\text{Prior}}$$

$$= \operatorname{argmax}_{c \in C} P(f_1, f_2, f_3, \dots, f_n | c) P(c)$$

$$C = \operatorname{argmax} \prod_{f \in F} P(c) P(f|c)$$

$$C = \operatorname{argmax} \log(P(c)) + \sum_{i \in \text{word positions}} P(f_i|c)$$

### Training on Naïve Bayes

- the goal is to learn probabilities

$$C = \operatorname{argmax} \log(P(c)) + \sum_{i \in \text{word positions}} P(f_i|c)$$

- $P(c)$  - what percentage of the documents in our training set are in each class C
  - o number of d in class c / number of documents (d)
- $P(f_i|c)$  - what fraction of times the word  $f_i$  appears among all words in all documents of class c
  - o ( count (  $f_i$ , c ) + 1 ) / sum(count(f, c) + |V| for the entire vocabulary (all words)

# Document classification

November 16, 2024 5:09 PM

Example IMDB ratings: response is the sentiment of the review (positive or negative)

- Find a way to featurize each document (movie review)
  - o define a set of predictors

Bag of words model

- Summarizes a document by the words present and ignores their context
- score each document for presence/absence of each of the words in a dictionary
- if the dict contains M words, then for each document we create a binary feature vector of length M and score a 1 for every word present
- \*\* limit dict i.e. to 10000 most frequently used words in the training set
  - o omit unknown words
- Build training set and test set balanced with regard to sentiment
  - o resulting training feature matrix X has dimensions # of reviews in the set x # of words in the dict
  - o only 1.3% of entries are non-zero -- sparse X

How to account for context

- bag of n grams model
  - o i.e. a bag of 2 grams records the consecutive co-occurrence of every distinct pair of words
    - blissfully long can be seen as a positive review
    - blissfully short can be seen as negative
  - o treat the document as a sequence, taking account all the words in the context of those that preceded and those that follow

# Naïve bayes

November 17, 2024 3:18 PM

Assume  $f_k$  is the density function for a multivariate normal random variable with class specific mean  $\mu_k$ , and shared covariance matrix  $\Sigma$

- These very strong assumptions allow us to simplify the problem of estimating K p-dimensional density functions ( $f_k$ 's) into estimating K p-dimensional mean vectors and ONE p x p-dimensional covariance matrix

## The Naïve bayes assumption

The Naïve bayes classifier takes a different tack for estimating  $f_1(x) \dots f_k(x)$

- Instead of assuming these functions belong to a particular family of distributions (i.e. multivariate normal)
- Make one assumption: Within the  $k^{\text{th}}$  class, the p predictors are independent
  - o for  $k = 1 \dots K$ 
$$f_k(x) = f_{k1}(x_1) \times f_{k2}(x_2) \times \dots \times f_{kp}(x_p)$$
where  $f_{kj}$  is the density function of the  $j^{\text{th}}$  predictor among observations in the  $k^{\text{th}}$  class
- By assuming that the p covariates are independent within each class, we don't have to worry about the association between the p predictors
  - o we assume there is no association between them
- Useful in settings where n is not large enough relative to p for us to effectively estimate the joint distribution of predictors within each class
- Introduces some bias, reduces variance

Posterior probability

$$\Pr(Y = k | X = x) = \frac{\pi_k \times f_{k1}(x_1) \times f_{k2}(x_2) \times \dots \times f_{kp}(x_p)}{\sum_{l=1}^K \pi_l \times f_{l1}(x_1) \times f_{l2}(x_2) \times \dots \times f_{lp}(x_p)}$$

for  $k = 1, \dots, K$ .

- to estimate the 1d density function  $f_{kj}$  using training data  $x_{1j}, \dots, x_{nj}$ , we can
  - o If  $X_j$  is quantitative, assume  $X_j | Y=k \sim N(\mu_{jk}, \sigma_{jk}^2)$ 
    - Assume that within each class, each predictor is drawn from a normal dist
    - assume predictors are independent
  - o If  $X_j$  is quantitative
    - non-parametric estimate for  $f_{kj}$
    - estimate  $f_{kj}(x_j)$  as the fraction of the training observations in the  $k^{\text{th}}$  class that belong to the same histogram bin as  $x_j$ 
      - kernel density estimator
  - o If  $X_j$  is qualitative
    - count the proportion of training observations for the  $j^{\text{th}}$  predictor corresponding to each class
    - i.e. if we have 100 observations of the  $k^{\text{th}}$  class, and the  $j^{\text{th}}$  predictor takes on values of 1, 2, and 3 in 32, 55, and 13 of those observations, then we can estimate

$$f_{kj}(x_j) = \begin{cases} 0.32 & \text{if } x_j = 1 \\ 0.55 & \text{if } x_j = 2 \\ 0.13 & \text{if } x_j = 3 \end{cases}$$

- Any classifier with a linear decision boundary is a special case of naïve bayes
  - LDA is a special case of naïve bayes
  - LSA assumes features are normally distributed with a common within-class covariance matrix
  - naïve bayes assumes independence of features
- KNN is better when the decision boundary is highly non linear
  - need a lot of observations relative to number of predictors
    - $n \gg p$
  - when decision boundary is non linear but  $n$  is modest or  $p$  is not very small, QDA may be better than KNN
  - KNN does not tell us which predictors are important

# L21 - Clustering

November 17, 2024 2:50 PM

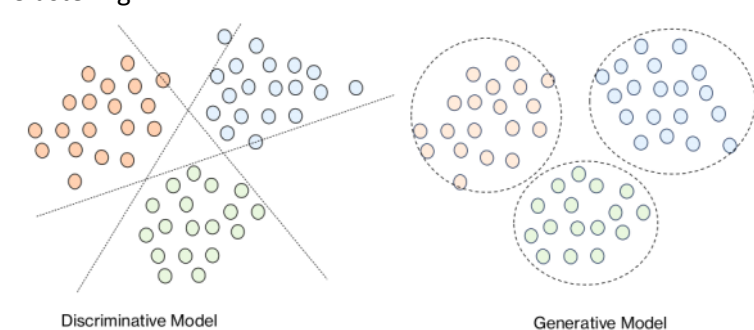
Text to vectors: Term Frequency - Inverse Document Frequency (TF - IDF)

- evaluate words wrt document (not position within a sentence)
- relative frequency of occurrence in a document
- can also account for number of documents in which each term can be found

Generative vs discriminative models

- Generative
  - o Can generate new data instances
  - o  $P(x, y) = X, Y \rightarrow [0, 1]$
  - o Capture the joint probability  $P(x, y)$
  - o Can work without labels, so can compute  $P(x)$
- Discriminative
  - o Discriminate between different kinds of data instances
  - o  $P(y|x) = X, Y \rightarrow [0, 1]$
  - o Capture conditional probability  $P(y \text{ given } x)$
  - o Needs labels to work

Clustering



- We may cluster the data into subgroups and partition the entire data distribution based on similarity
- The goal of clustering is to find homogenous subgroups among observations
- When a new data instance arrives, we find the subgroup it may belong to or is closest to
- Relies on data distribution to build a classification model
- They are challenging, tackle difficult tasks

Types of clustering

- K-means
  - o Let  $C_1, C_2, \dots, C_k$  be sets containing observations in each cluster
  - o These sets must satisfy...
    - Each observation must belong to at least one of the  $k$  clusters
    - The clusters are non overlapping, and no observation belongs to more than one cluster
  - o The variation of a cluster is  $W(C_k)$
  - o Good clustering is one for which the within-cluster variation is minimal
    - Optimize:
$$\text{minimize } \left\{ \sum_{k=1}^K W(C_k) \right\}$$
$$W(C_k) = \frac{1}{|C_k|} \sum_{i, t \in C_k} \sum_{j=1}^p (x_{ij} - x_{tj})^2$$
      - Where  $|C_k|$  is the number of observations in the  $k^{\text{th}}$  cluster
  - o Defining similarity

- Feature similarity using Euclidian distance
  - distance between 2 pts  $d(p, q) = |p, q|$
  - In high dimensionality  $d(p, q) = \sqrt{(p - q)^2}$
- Training k means
  1. Randomly assign a number, from 1 to K, to each observation. Serve as initial cluster assignments
  2. Iterate until the cluster assignments stop changing
    - a) for each of the K clusters, find the cluster centroid
      - ◆ the  $k^{\text{th}}$  cluster centroid is the vector of the p feature means for the observations in the  $k^{\text{th}}$  cluster
$$\mu_i = \frac{1}{|S_i|} \sum_{x \in S_i} x$$
    - b) Assign each observation to the cluster whose centroid is the closest
      - ◆ (Euclidian distance is smallest)



# K means clustering

November 17, 2024 3:56 PM

- For partitioning the dataset into K distinct, non-overlapping clusters
- let  $c_1 \dots c_k$  denote sets containing the indices of the observations in each cluster
  1.  $C_1 \cup C_2 \cup \dots \cup C_k = \{1, \dots, n\}$   
Each observation belongs to at least one of the K clusters
  2.  $C_k \cap C_{k'} = \emptyset$  for all  $k \neq k'$   
The clusters are non-overlapping: no observation belongs to more than one cluster
- A good cluster = minimal within-cluster variation
  - amount by which observations within a cluster differ from each other
  - Partition observations into K clusters s.t. sum of within-cluster variations across K clusters is minimal
    - calculate variation using euclidian distance
- $$\text{minimize}_{C_1, \dots, C_K} \left\{ \sum_{k=1}^K \frac{1}{|C_k|} \sum_{i, i' \in C_k} \sum_{j=1}^p (x_{ij} - x_{i'j})^2 \right\}$$
- K means clustering algorithm
  1. Randomly assign a number from 1 to K, to each observation
    - Initial cluster assignments
  2. Iterate until the cluster assignments stop changing
    - a. For each cluster, compute the centroid
      - centroid = vector of the p feature means for the observations in the kth cluster
    - b. Assign each observation to the cluster whose centroid is closest
- $$\frac{1}{|C_k|} \sum_{i, i' \in C_k} \sum_{j=1}^p (x_{ij} - x_{i'j})^2 = 2 \sum_{i \in C_k} \sum_{j=1}^p (x_{ij} - \bar{x}_{kj})^2,$$
  - where  $\bar{x}_{kj} = \frac{1}{|C_k|} \sum_{i \in C_k} x_{ij}$
- As the algorithm is run, reallocating the observations can only improve the clustering
  - the clustering obtained will be continuously improved until the result no longer changes
  - finds a local min and stays there
    - because it does not necessarily find a global min, depends on the initial cluster assignment
  - run the algo multiple times with different init configs to select the best solution

# L22 - Clustering and dimensionality reduction

November 17, 2024

4:38 PM

## Types of clustering

- K-means
- Hierarchical
  - o Uses tree based representation of observations called Dendrograms
  - o The lower in the tree, fusions occur, the more similar the groups of observations are to each other
  - o Most common type is bottom-up (agglomerative)
  - o For n observations,  $2^{n-1}$  reordering of dendrograms

## Hierarchical clustering algorithm

1. Begin with n observations and a measure (i.e. Euclidian distance) of all the  $\binom{n}{2} = \frac{n(n-1)}{2}$  pairwise dissimilarities
  - a. treat each observation as its own cluster
2. For i = n, n-1, ..., 2:
  - a. Examine all pairwise intercluster dissimilarities among the i clusters
    - i. identify the pair that are the most similar (least dissimilar)
    - ii. fuse the two similar clusters
    - iii. the dissimilarity between those two clusters indicates the height in the dendrogram at which the fusion should be placed
  - b. compute the new pairwise intercluster dissimilarities along the i-1 remaining clusters

## Design considerations

- Linkage techniques
  - o Complete: maximal intercluster dissimilarity
  - o Single: minimal intercluster dissimilarity
  - o Average: mean intercluster dissimilarity
  - o Centroid: dissimilarity between the centroid
- Measures
  - o Euclidian distance
  - o correlation based measures

## Principal component analysis

- Suppose we wish to visualize n observations with measurements on a set of p features
- not all of p dimensions are equally interesting
- PCA find a low dimensional representation of the data set that contains as much as possible of the variation
  - o seeks a small number of dimensions that are as interesting as possible
- First principal component
  - o Normalized linear combination of the features that has the largest variance
$$Z = \phi_{11}X_1 + \phi_{21}X_2 + \dots + \phi_{p1}X_p$$
$$\phi_{11} = \text{loadings of the PCA}$$
    - define a direction in feature space along which the data varies the most
  - o Normalized linear combination of the features
$$\sum_{j=1}^p \phi_{j1}^2 = 1$$

# Hierarchical clustering

November 17, 2024 4:50 PM

K-clustering requires us to pre-specify the number of clusters K

## Hierarchical clustering

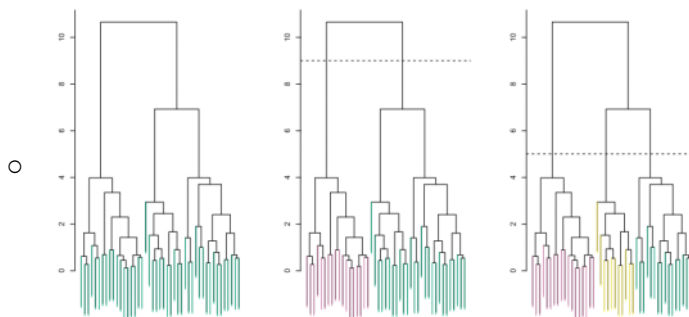
- an alternative approach which does not require us to commit to a particular choice K
- Results in a tree-based representation called a dendrogram
- Bottom-up clustering
  - o most common type
  - o dendrogram is built from the leaves up, combining clusters up to the trunk

## Interpreting the dendrogram

- each leaf represents an observation
- as we move up in the tree, some leaves *fuse* into branches
  - o correspond to similar observations
- as we move more up the tree, some branches fuse
- these fused branches fuse and so on and so on
- the earlier the fusions occur, the more similar the groups of observations are
- observations that fuse later can be very different
- There are  $2^{n-1}$  possible reorderings of the dendrogram, where n is the number of leaves
  - o at the n-1 point where fusion occurs, the positions of two fused branches can be swapped without changing the meaning of the dendrogram
- do not draw conclusions about similarity of two observations based on their proximity along the horizontal axis
  - o look at vertical axis position where branches are fused

## Clustering with a dendrogram

- Cutting the dendrogram at various heights results in clusters



- number of observations beneath the cut = cluster
- can cut the dendrogram to get any number of clusters between 1 and n
- Sometimes, data cannot be well represented by hierarchical clustering
  - o i.e. the best division into 3 groups is not the result from taking the best division into 2 groups and splitting up 1 of them into 2
    - i.e. suppose observations correspond to group of men and women evenly split among americans, japanese, and french
      - best division into 2 is by gender
      - best division into 3 is by nationality

## The hierarchical clustering algorithm

- Simple algo
  - measure dissimilarity (using euclidian distance)
1. each observation is treated as its own cluster
  2. the two most similar clusters are fused so that there are now n-1 clusters
  3. The next two clusters are fused again, n-2 clusters

4. Repeat until all of the observations belong to one single cluster
  - How do we define dissimilarity between two clusters if one or both contains multiple observations
    - o linkage: defines dissimilarity between two groups of observations
      - Complete, average, single, and centroid
      - Average and complete are preferred

#### Measuring dissimilarity

- alternatives to euclidian distance
- correlation-based distance
  - o considers two observations to be similar if their features are highly correlated
  - o may be far apart wrt euclidian distance
  - o focuses on the shape of observations
- choose the measurement type which fits the problem well
- i.e. shopping preferences - assume data is a matrix where rows are shoppers and columns are items available for purchase, entries is how many times they bought each item
  - o if euclidian distance is used, shoppers who don't buy many things will be clustered together
  - o if correlation-based distance is used, shoppers with similar preferences will be clustered

#### Issues with clustering

- Small decisions with big consequences
  - o Should the observations or features first be standardized in some way?
    - i.e. scaling
  - o What dissimilarity measure should we use?
  - o What type of linkage should be use?
  - o Where do we cut the dendrogram?
  - o In k-means clustering, how many clusters should we look for in the data?
  - o In practice, try several different choices and observe results - there is no single right answer
- Validating the obtained clusters
  - o do the clusters represent the true subgroups of the data
    - could be a result of clustering the noise
  - o no single best approach on validaitng this
- Other considerations
  - o both k means and hierarchical will assign each observation to a cluster
    - sometimes this may not be appropriate
    - i.e. most of the observations truly belong to a small number of unknown subgroups
    - a small subset of the observations are very different from eachother and from all other observations
    - since every observation is forced into a cluster, the clusters found might be distorted due to the presence of outliers that don't belong to any cluster
    - mixture approaches address this, i.e. soft k means clustering
  - o Clustering methods are not robust
    - cluster n observations, remove a subset, cluster again: you'll likely end up with different clusters

# L23 - Unsupervised Learning: PCA

November 17, 2024 5:43 PM

## Principal component analysis

- reduction of dimensionality in visualization of features
- Suppose we wish to visualize n observations with measurements on a set of p features
- not all of p dimensions are equally interesting
- PCA find a low dimensional representation of the data set that contains as much as possible of the variation
  - o seeks a small number of dimensions that are as interesting as possible

### - First principal component

- o Normalized linear combination of the features that has the largest variance

$$Z = \phi_{11}X_1 + \phi_{21}X_2 + \dots + \phi_{p1}X_p$$

$\phi_{11}$  = loadings of the PCA

- define a direction in feature space along which the data varies the most

- o Normalized linear combination of the features

$$\sum_{j=1}^p \phi_{j1}^2 = 1$$

$$Z_m = \sum_{j=1}^p \phi_{jm}X_j \rightarrow y_i = \theta + \sum_{m=1}^M \theta_m Z_m + \epsilon_i$$

can solve this with OLS

- o Solve the optimization problem:

$$\underset{\phi_{11}, \dots, \phi_{p1}}{\text{maximize}} \left\{ \frac{1}{n} \sum_{i=1}^n \underbrace{\left( \sum_{j=1}^p \phi_{j1} x_{ij} \right)^2}_{Z_{i1}^2} \right\} \text{ subject to } \sum_{j=1}^p \phi_{j1}^2 = 1$$

### - Second principal component

- o normalized linear combination of the features that has the largest variance out of all linear combinations that are uncorrelated with  $Z_1$

$$Z = \phi_{12}X_1 + \phi_{22}X_2 + \dots + \phi_{p2}X_p$$

$\phi_{i2}$  = loadings of the PCA

- o Constraining  $Z_2$  to be uncorrelated with  $Z_1$  is equivalent to constraining the direction of  $\phi_1$  to be orthogonal to the direction of  $\phi_2$

$$\underset{\phi_{11}, \dots, \phi_{p1}}{\text{maximize}} \left\{ \frac{1}{n} \sum_{i=1}^n \left( \sum_{j=1}^p \phi_{j2} x_{ij} \right)^2 \right\} \text{ subject to } \sum_{j=1}^p \phi_{j2}^2 = 1$$

## PCA

- For multidimensional data, there are multiple principal components
- Principal components can be used to produce low-dimensional views
  - o It is not feature selection, we consider linear combination of all features
- Eigen decomposition method is used to solve PCA
  - o PC direction are computed using the order sequence of eigenvectors of the matrix  $X^T X$
- For a set of p features and n observations, there are at most  $\min(n-1, p)$  principal components
  - o we use the smallest subset, determined by proportion of variance

- Before PCA is performed, the variables should be centered to have mean=0
  - o PCA results are heavily dependent on scaling

# Dimensionality reduction

November 17, 2024 6:10 PM

Instead of taking a subset or shrinking coefficients, we can transform the predictors and fit a OLS model to reduce dimensions

Let  $Z_1, Z_2, \dots, Z_M$  represent  $M < p$  linear combinations of our original  $p$  predictors. That is,

$$Z_m = \sum_{j=1}^p \phi_{jm} X_j \quad (6.16)$$

for some constants  $\phi_{1m}, \phi_{2m}, \dots, \phi_{pm}$ ,  $m = 1, \dots, M$ . We can then fit the linear regression model

$$y_i = \theta_0 + \sum_{m=1}^M \theta_m z_{im} + \epsilon_i, \quad i = 1, \dots, n, \quad (6.17)$$

- the thetas are the regression coefficients

## Dimension reduction

- reduces the problem of estimating the  $p+1$  coefficients (betas) to the simpler problem of estimating the  $M+1$  coefficients (thetas), where  $M < p$ 
  - o the dimensions of the problem have been reduced from  $p+1$  to  $M+1$
- Dimension reduction serves to constrain the estimated  $\beta_j$  coefficients since now they must take their form:

$$\beta_j = \sum_{m=1}^M \theta_m \phi_{jm}.$$

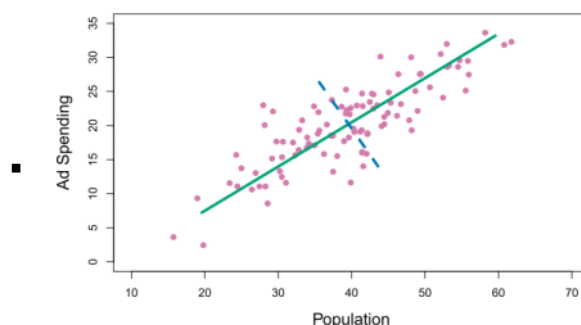
- o potentially increased Bias but significantly reduced variance
- If  $M = p$  and all  $Z_m$  are linearly independent, then this poses no constraints and no dimension reduction occurs

## How does dimension reduction work

1. The transformed predictors  $Z_1, Z_2, \dots, Z_M$  are obtained
2. the model is fit using these  $M$  predictors
  - the choice of  $Z_1, Z_2, \dots, Z_M$  (or the phi's) can be achieved in different ways
    - o principal components and partial least squares

## Principal components regression/PCA

- reducing the dimensions of an  $n \times p$  matrix  $X$
- First principal component direction of the data
  - o the direction along which the observations vary the most
  - o the line s.t. if we project the points onto the line, the points will have the largest possible variance
    - finding the location on the line which is closest to the point



- fpc = green line
- o First principal component eqn

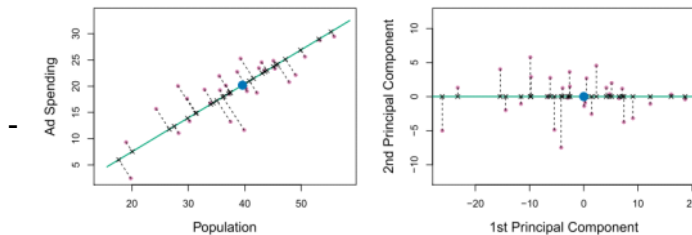
$$Z_1 = \phi_{11} * (pop - mean_{pop}) + \phi_{21} * (ad - mean_{ad})$$

where  $\phi_{11} = 0.839$  and  $\phi_{21} = 0.544$

- The idea: out of every possible linear combination of pop and ad s.t.  $\phi_{11}^2 + \phi_{21}^2 = 1$ , this linear comb. yields the highest variance
  - o maximized
$$\text{Var}(\phi_{11} \times (\text{pop} - \overline{\text{pop}}) + \phi_{21} \times (\text{ad} - \overline{\text{ad}}))$$
- Principal component scores
- $z_{i1} = \phi_{11} * (pop_i - mean_{pop}) + \phi_{21} * (ad_i - mean_{ad})$

#### Alternate interpretation

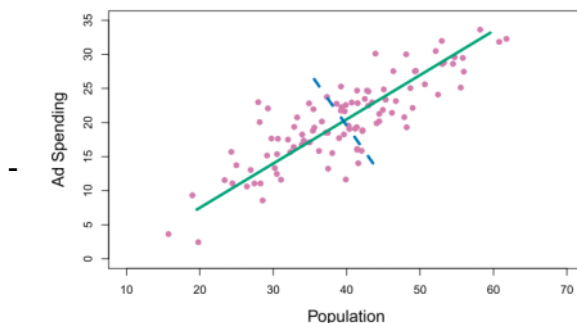
- the first principal component vector is the line that is as close as possible to the data
  - o minimizes the sum of the squared perpendicular distances between each pt and the line



- the first principal component appears to capture most of the info contained in the pop and ad predictors
  - o the plots show a strong relationship between the first principal component and the two features
  - o pop and ad have an approximately linear relationship

#### The second principal components

- we can have up to p distinct principal components
- $Z_2$  - the 2nd principal component: a linear combination of the variables that is uncorrelated with  $Z_1$  and has the largest variance subject to this constraint



- spc = blue line
  - must be perpendicular to the fpc direction
  - Second principal component formula
- $$Z_2 = 0.544 * (pop - mean_{pop}) - 0.839 * (ad - mean_{ad})$$
- since the data has 2 predictors, fpc and spc contain all the information that is in pop and ad
    - o the fpc will contain the most information

#### Principal components regression

- construct the first M principal components  $Z_1 \dots Z_m$  and use these components as predictors in a linear regression model that is fit using least squares
- often a small number of principal components suffice to explain most of the variability in the data
- We assume that the directions in which  $X_1, \dots, X_p$  show the most variation are the directions that are associated with Y
  - o approximation
- Fitting the least squares model to  $Z_1 \dots Z_m$  will lead to better results than fitting to  $X_1 \dots X_p$ , since most or all of the information in the data pertaining to Y is contained in the Z's



- by estimating fewer coefficients, we can mitigate overfitting ( $M \ll p$ )
- More principal components = less bias but more variance
  - MSE still a U shape
- PCR does best when the first few principal components are sufficient to capture most of the variation in the predictors as well as the relationship with the response
- NOT a feature selection method

# L24 - Visualization

November 20, 2024 3:34 PM

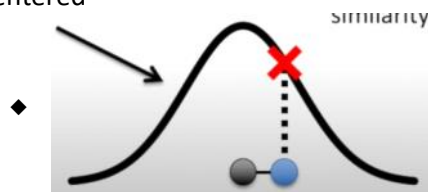
- Need visualization for
  - o understanding training data
  - o inspecting model
  - o communicating model results
  - o dealing with high dimensional data

## High-dimensionality: T-SNE

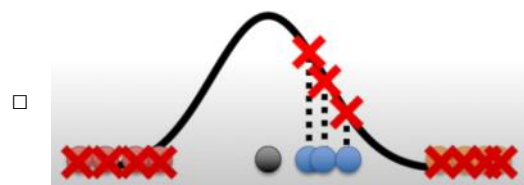
- Gaussian dist
  - o a normal (gaussian) distribution is a type of continuous probability distribution for a real-valued random variable
- Kullback - Leibler divergence
  - o a measure of how one reference probability dist P is different from a second probability dist Q
  - o Also called relative entropy and I-divergence

$$D_{KL}(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \log \left( \frac{P(x)}{Q(x)} \right)$$

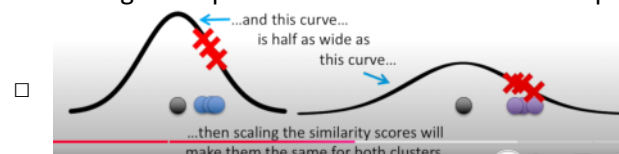
- T-sne
  - o finds a way to project the data into a low dimensional space while preserving clustering
    - i.e. 2d graph onto 1d line
      - start by randomly placing all points onto the line
      - tsne will move these points a little at a time until they have been clustered
        - ◆ points are attracted to other points of the same cluster
        - ◆ points are repelled by points of distant/different clusters
        - ◆ think one point at a time
        - ◆ attraction is stronger than repelling
    - determining the similarity between points
      - measure the distance between 2 points
      - plot that distance on a normal/gaussian curve, with the point of interest centered

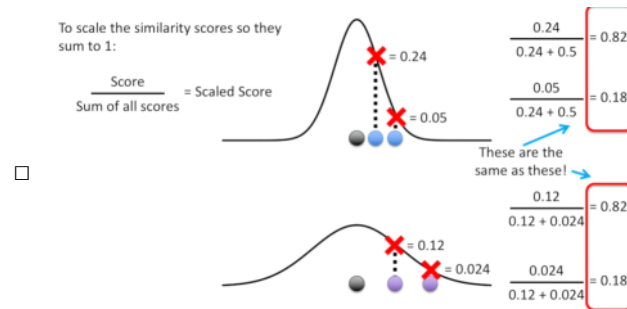


- x is the unscaled similarity between the point of interest (black) and the other point (blue)
- calculate unscaled similarity for black point and all other points one by one

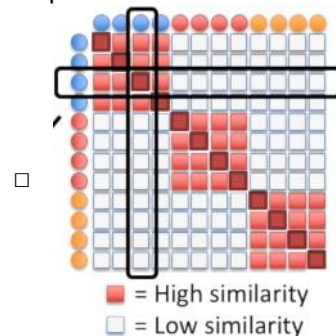


- rescale all unscaled similarities so they add to 1
  - width of normal curve depends on density of the cluster
  - scaling will equalize distances for dense and sparse clusters

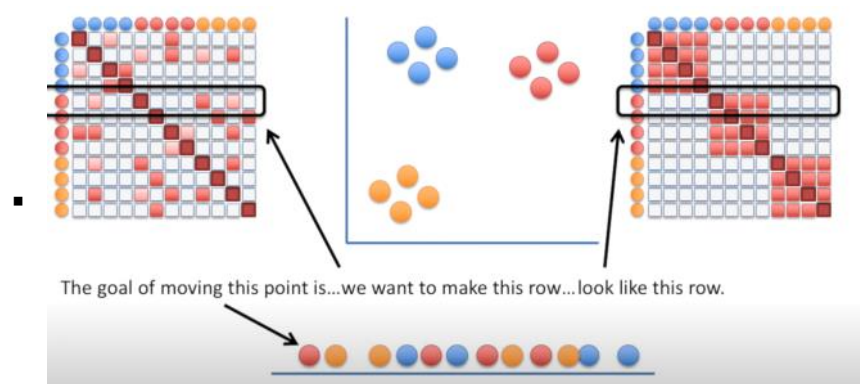
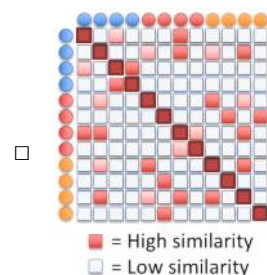




- ◆ perplexity: expected density
- because normal dist is based on density of nearby points (cluster), similarity of point A1 to B1 might be different than similarity of point B1 to A1 because cluster A and cluster B may have different densities
  - ◆ average similarities from both directions to fix this
- end up with a matrix of similarity scores



- Now, go back to the line example
  - For all the points on the line, measure the distance as described above using a t-distribution instead of a normal dist



- move the points a little at a time, and at each step it chooses a direction that makes the t-dist matrix look more like the n-dist matrix

- for n high dimensionality points, stochastic neighbour embedding (sne) starts converting the high dimensional euclidian distances between datapoints into conditional probabilities that represent similarities
- the similarity of  $x_j$  to  $x_i$  is  $p_{j|i}$ 
  - is the probability that  $x_i$  would pick  $x_j$  as its neighbour if neighbours were picked in proportion to their probability density under a gaussian centered at  $x_i$
- For nearby points,  $p_{j|i}$  is relatively high
- for widely separated points,  $p_{j|i}$  is extremely small

$$p_{j|i} = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|x_i - x_k\|^2 / 2\sigma_i^2)}$$

#### - Process

- since we are only interested in modeling pairwise similarities,  $p_{i|i}$  is set to 0
- consider the low dimensionality counterparts  $(y_i, y_j)$  of the high dimensionality points  $(x_i, x_j)$
- The similar the conditional probability  $q_{j|i}$  can be computed
- Gaussian variance  $\sigma$  is set to  $\sqrt{1/2}$

$$p_{j|i} = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|x_i - x_k\|^2 / 2\sigma_i^2)} \quad q_{j|i} = \frac{\exp(-\|y_i - y_j\|^2)}{\sum_{k \neq i} \exp(-\|y_i - y_k\|^2)}$$

#### - How it works

- SNE minimizes the sum of divergences over all datapoints using a gradient descent method where...
  - $P_i$  represents cond prob dist over all other datapoints given data point  $x$
  - $Q_i$  represents cond prob dist over all other map points given map point  $y_i$

$$C = \sum_i KL(P_i || Q_i) = \sum_i \sum_j p_{j|i} \log \frac{p_{j|i}}{q_{j|i}} \quad \frac{\delta C}{\delta y_i} = 2 \sum_j (p_{j|i} - q_{j|i} + p_{i|j} - q_{i|j})(y_i - y_j)$$

Gradient updates  $\gamma^{(t)} = \gamma^{(t-1)} + \eta \frac{\delta C}{\delta \gamma} + \alpha(t) (\gamma^{(t-1)} - \gamma^{(t-2)})$  Without symmetric adjustments

$$C = \sum_i KL(P_i || Q_i) = \sum_i \sum_j p_{j|i} \log \frac{p_{j|i}}{q_{j|i}} \quad \frac{\delta C}{\delta y_i} = 4 \sum_j (p_{ij} - q_{ij})(y_i - y_j)$$

Gradient updates  $\gamma^{(t)} = \gamma^{(t-1)} + \eta \frac{\delta C}{\delta \gamma} + \alpha(t) (\gamma^{(t-1)} - \gamma^{(t-2)})$  With symmetric adjustments

- Parameters to be selected
  - variance of the gaussian that is centered over each high-dim datapoint
- SNE performs binary search for the value  $\sigma_i$  that produces a  $P_i$  with a fixed perplexity that is specified by the user
  - $Perp(P_i) = 2^{H(P_i)}$
- Where shannon entropy of  $P_i$  measured in bits
  - $H(P_i) = - \sum_j p_{j|i} \log_2 p_{j|i}$
- Step size for gradient descent

# !L25 - ML OPs

November 20, 2024

3:35 PM

# L26 - Intro to reinforcement learning

November 23, 2024 1:28 PM

- A software agent
  - o makes observations
  - o takes actions within an environment
  - o receives a reward
    - the reward may be positive - award for taking the right action
    - the reward may be negative - punishment for taking the wrong action
- Agent must maximize the reward and minimize punishment
- entirely trial and error
- Pac man
  - o environment is the simulation of the game
  - o agent controls pac man's movements
  - o actions are possible directions in which the movements can occur
    - up left right down
  - o reward is points earned/game score
  - o observation is screenshots of the game stages

## Policy

- RL needs policy function to determine behaviour of the agent
- policy defines agent's way of behaving at a given time
- mapping from perceived states (environment) to actions the agent must take when that state is encountered
- can be a simple hash table or a complex search algorithm

## Reward

- at each time step, agent receives a numeric reward and wants to maximize it
- reward depends on current action and current observation
- agent may learn a policy based on reward
  - o not learn reward based on policy
- if in the current state, the policy selected gives low reward, the policy is changed for the future state
- reward signal indicates what is good in current state

## Value

- indicates what is advantageous in the long run
- long-term-desirability of a given state in the env
- a state with low immediate reward does not necessarily yield a bad score, maybe it is helpful in the long run
- value function is estimated through observations and reward
- goal is to seek state of highest value, not highest reward

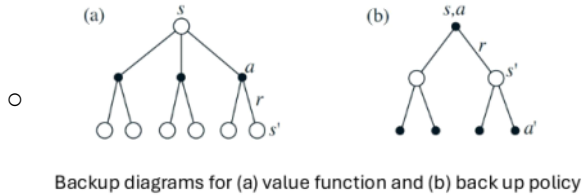
## Model

- simulation of the environment and mimic its behaviour in real time
- for given state and action, model might predict the next state and reward
- models are useful for planning a set of actions to anticipate states an agent might experience even before they are encountered
- model based methods use a model and planning
  - o model free methods do not

## Learning

- the idea is to play many games

- we observe the result from each move many times, select states that lead to the greatest value
- as we progress, we change the values of our current state
- to make more accurate estimates of the probability of success, we back up the value of state after each greedy move



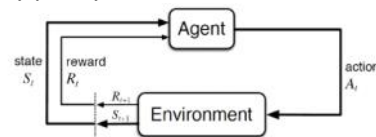
- $s$  = state before greedy move,
- $s'$  = state after the move,
- $V$  = estimated value of  $s$   $V(s) \leftarrow V(s) + \alpha [V(s') - V(s)]$

- if the step size parameter is reduced properly over time, this method converges, for any fixed opponent
- if the step size parameter is not reduced all the way to zero over time
  - player also plays well against opponents that slowly change their way of playing

### Policy evaluation

- how to we compute state-value of an arbitrary policy

$$\begin{aligned}
 v_{\pi}(s) &= \mathbb{E}_{\pi}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \mid S_t = s] \\
 &= \mathbb{E}_{\pi}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s] \\
 &= \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma v_{\pi}(s')],
 \end{aligned}$$



- if the env's dynamics are completely known, then the above is a system of  $|S|$  simultaneous linear equations
- Use the bellman equation
  - averages over all the possibilities, weighting each by its probability of occurring

$$\begin{aligned}
 &= \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma v_{\pi}(s')]
 \end{aligned}$$

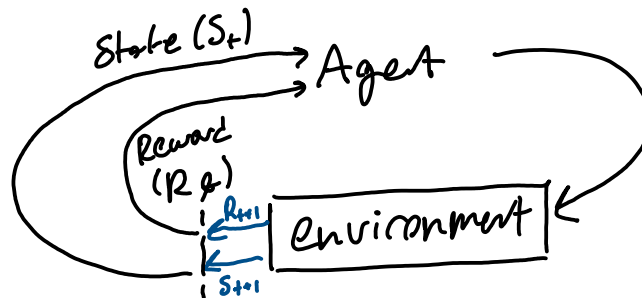
$$\begin{aligned}
 v_{k+1}(s) &= \mathbb{E}_{\pi}[R_{t+1} + \gamma v_k(S_{t+1}) \mid S_t = s] \\
 &= \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma v_k(s')]
 \end{aligned}$$

# L28 - Reinforcement

November 20, 2024 3:36 PM

## Markov decision process

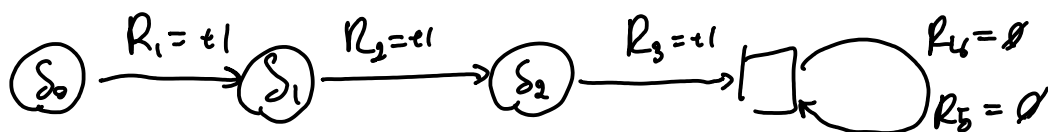
- Reward received now versus in the future does not carry the same weight



$0 < \gamma < 1$  is the discount rate

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{k-1} R_{t+k}$$

- a state signal succeeds in retaining all relevant info is said to be Markov, or to have the Markov property
  - o Markov signal: retains all relevant info



- i.e. a game of checkers, screenshot of pac man, can see all important information at once
- When is it not a markov decision process?
  - o The marble draw: when there is a connection between the previous marble draw and the current marble draw, this is not a markov process
- The environment's response at  $t + 1$  depends only on the state and action representations at  $t$
- The markov property is important because decisions and values are assumed to be a function of ONLY THE CURRENT STATE

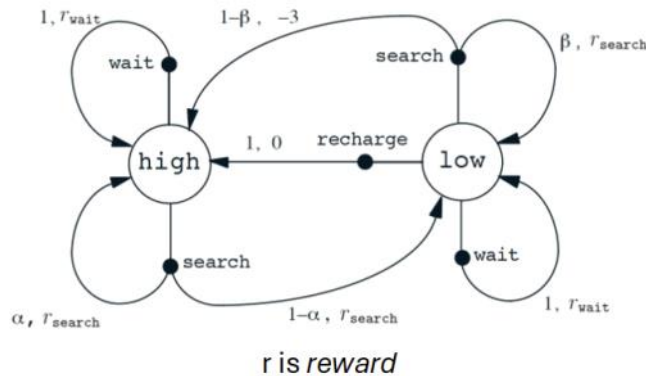
## Markov decision process

- a reinforcement learning task that satisfies the markov property (MDP)
- i.e. a robot throwing out trash
  - o Energy level = {high, low} (battery)
  - o A {high} = {search, wait} where search and wait are actions
  - o A {low} = {search, wait, recharge}
  - o Can use this to make a hash table
 
$$p(s', r | s, a) = \Pr\{R_{t+1} = r, S_{t+1} = s' \mid S_t = s, A_t = a\}$$



$s$	$s'$	$a$	$p(s' s, a)$	$r(s, a, s')$
high	high	search	$\alpha$	$r_{\text{search}}$
high	low	search	$1 - \alpha$	$r_{\text{search}}$
low	high	search	$1 - \beta$	$-3$
low	low	search	$\beta$	$r_{\text{search}}$
high	high	wait	1	$r_{\text{wait}}$
high	low	wait	0	$r_{\text{wait}}$
low	high	wait	0	$r_{\text{wait}}$
low	low	wait	1	$r_{\text{wait}}$
low	high	recharge	1	0
low	low	recharge	0	0

- Can use to make a flow diagram



### Value functions and Dynamic programming

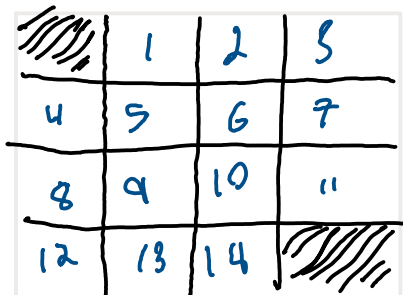
- value functions are estimated for states (or state-action pairs) to estimate how good it is for an agent to be in each state
  - or how good it is to perform an action while in a state
- State value function and action value function is optimized based on a certain policy and is estimated from experience

- state-value function:  $V_{\pi}(s)$   
- action-value function:  $q_{\pi}(s)$   
- policy:  $\pi$

- Bellman equation:  $V_{\pi}(s) = \sum_a \pi(s|a) \sum_{s', r} p(s', r|s, a) [r + \gamma V_{\pi}(s')]$

### Policy iteration

- nonterminal states  $\{1, 2, \dots, 14\}$   $R = -1$  on all transitions until the terminal state is reached
- actions: move up down left right
- we need a policy and a value function



$p(6|5, \text{right}) = 1$   
 $p(10|5, \text{right}) = 0$   
 $p(7|7, \text{right}) = 1$   
 $r(s, a, s') = 1$

- Start with random policy, then create greedy policy wrt random policy

$V_k$  for the Random Policy

Greedy Policy w.r.t.  $V_k$

○

0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0

	↕	↕	↕
↕	↕	↕	↕
↕	↕	↕	↕
↕	↕	↕	

- Iteration

input  $\pi$ , the policy to be evaluated

init an array  $V(s) = 0$ , for all  $s \in S$

Repeat

- $\Delta \leftarrow 0$
- For each  $s \in S$ :
  - $v \leftarrow V(s)$
  - $V(s) \leftarrow \sum_a \pi(s|a) \sum_{s',r} p(s',r|s,a)[r + \gamma v_\pi(s')]$
  - $\Delta \leftarrow \max(\Delta, |v - V(s)|)$

Until  $\Delta < \theta$  (a small positive number)

Output  $V \approx v_\pi$

$$\pi_0 \xrightarrow{E} v_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} v_{\pi_1} \xrightarrow{I} \pi_2 \xrightarrow{E} \dots \xrightarrow{I} \pi_* \xrightarrow{E} v_*$$

Limitations of policy iteration

- each iteration involved policy evaluation
  - heavy computation

Value iteration

- Initialize an array  $V(s) = 0$ , for all  $s \in S$
- Repeat
  - $\Delta \leftarrow 0$
  - For each  $s \in S$ :
    - $v \leftarrow V(s)$
    - $V(s) \leftarrow \max_a \sum_{s',r} p(s',r|s,a)[r + \gamma v_\pi(s')]$
    - $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
- Until  $\Delta < \theta$  (a small positive number)
- Output a deterministic policy,  $\pi$ , such that

$$\pi(s) = \operatorname{argmax}_a \sum_{s',r} p(s',r|s,a)[r + \gamma v_\pi(s')]$$

# L29 - RNN

November 23, 2024

2:58 PM

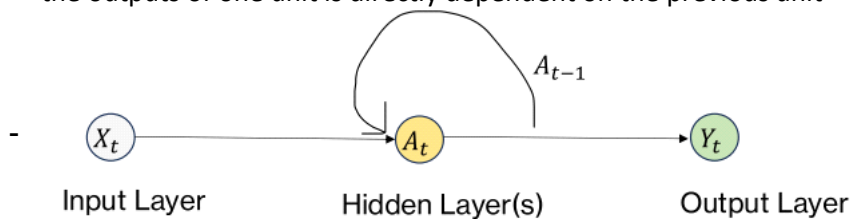
- when sequence is important in datasets
- when applying ML to sequences
  - o input variable is a sequence
  - o target var is also a sequence
- Rarely do we have fixed target sequences
  - o in this case, we predict the next item in the sequence at next time step  $t+1$  based on input data until time step 1

## Feedforward NN

- fixed size input vectors with associated weights to capture all relevant aspects of an example at once

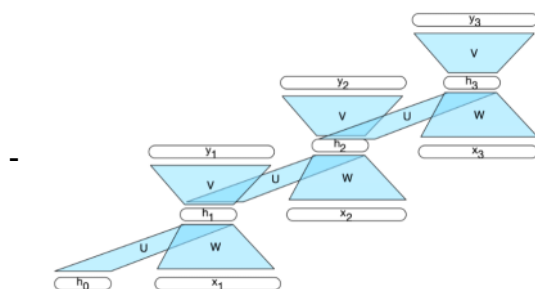
## RNN

- recurrent neural network contains a feedback connection within its network connections
- the outputs of one unit is directly dependent on the previous unit

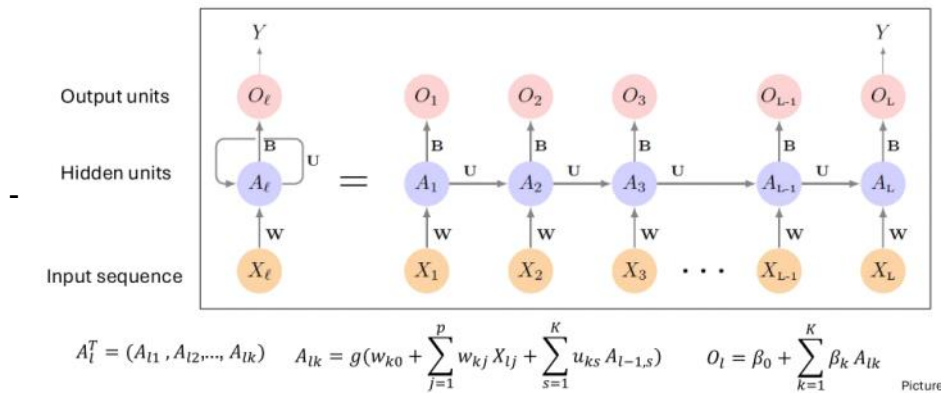


## Inference

- the mapping between inputs and outputs is almost like traditional NNs
- $h_t = g(Uh_{t-1} + Wx_t)$
- $y_t = \text{softmax}(Vh_t)$
- Given a hidden state in our model, allows storage of info
- goal is to get a good estimate of prob dist over the space of hidden state vectors to compute the next output



- Giving a hidden state to our model allows for storage of info
- get a good enough estimate of prob dist over the space of hidden state vectors, to compute the next output
  - $h_0 \leftarrow 0$ 
    - for  $i \leftarrow 1$  to length( $x$ ) do
      - $h_i \leftarrow g(Uh_{i-1} + Wx_i)$
      - $y_i \leftarrow f(Vh_i)$
    - return  $y$



### Text classification

- use embedding spaces instead of onehot
- embedding space can be learned or weights can be static
  - o popular frozen weights: word2vec, GloVe
    - GloVe: unsupervised learning algo for obtaining vector representations for words

```
word_to_ix = {"how": 0, "are": 1, "you": 2}
embeds = nn.Embedding(3, 5) # 3 words in vocab, 5 dimensional embeddings
lookup_tensor = torch.tensor([word_to_ix["how"]], dtype=torch.long)
how_embed = embeds(lookup_tensor)
print(how_embed)
```

### Time-series forecasting

- predicting trading trend/stocks to buy/stock price is very hard
- predicting trade volume is more reasonable

Read chapter 10 of textbook