

9/22

## Today

- CS AG update
  - Homework? GitHub
  - Review from Last year
    - Prob/stats
    - NN
    - NLP
- 

## NLP

- Translation
- Summarization
- Sentiment Analysis

### ① Representations of Language / Feature Extraction

\* The movie was great + , 1 Label

→ POS tagging the - article  
movie - noun  $\rightarrow$  NN

{ great, good ... }  $\leftarrow$  +

{ bad, terrible ... }  $\leftarrow$  -

The movie was not great - + ?

Bag of Words :  $\bar{x}$  : [The movie was great]

Corpus :  $\mathcal{V} \leftarrow f(\bar{x})$  : [  $\begin{matrix} \text{the} & \text{movie} & \text{spectacular} & \text{great} \\ 1 & 1 & 0 & 1 \end{matrix}$  ]

$|V| \times d$  vectors

- Sparse vectors

- Dense vectors

↳ - Counter

'the' : 1

'movie' : 1

⋮

great	great
bad	<u>GREAT</u>

## Preprocessing

① tokenizing → sentence → tokens  
great great! ←

② stopwords : a, the X

③ lowercasing : great GREAT  
↳ ✓

④ rare words : counter ← 10,000 words

↳ lowercasing :

The movie was GREAT → the movie was great

---

nlTK

↳ - Domain-specific tasks

$f(\bar{x}) \rightarrow$  dense feature vector

$\bar{w} \rightarrow$  weight vector

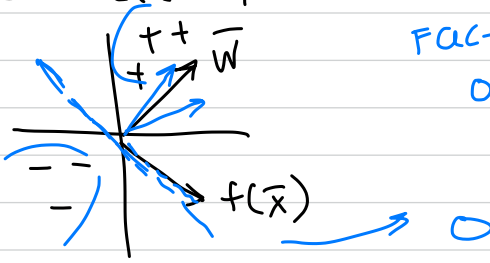
$$\bar{w} \cdot \underline{f(\bar{x})}$$

$$> 0 \Rightarrow +$$

$$< 0 \Rightarrow -$$

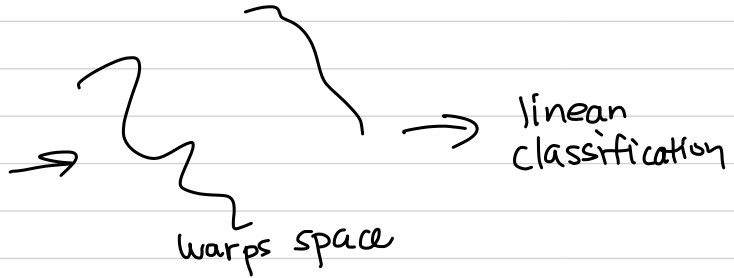
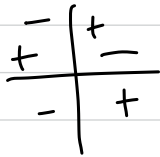
$$\begin{bmatrix} 4 \\ 3 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 2 \end{bmatrix} = 4 + 6 = 10$$

Linear classification



Fact: dot product of orthogonal vectors is 0

Neural Nets



---

- (really, good)  $\rightarrow$  Counter (not, good)  $\rightarrow$  : 1  
- (not, good)

The movie was not good

Bigrams

single words:  
Unigrams

LR

- Unigram 81.7% Pang et al. (2002)
- Bigrams 77.3%
- Unigrams + Bigrams 80.6%

Counter

- "the" : 1
- "(the, movie)" : 1

Kim (2014)

CNN

$\sim 83\%$

---

update weights in training

- SGD : stochastic gradient descent

for  $t$  in range  $(0, \text{epochs})$ :

for  $i$  in range  $(0, D) \leftarrow D$  is # of labeled

sample  $j \sim \{0, 1, \dots, D-1\}$  data points

- predict

$$\bar{w} \leftarrow \bar{w} - \alpha \cdot \frac{\partial \text{loss}}{\partial \bar{w}}$$

↑  
step size / learning rate

# Probability / Statistics Review

## Definitions

(R.V.) Random Variable ( $X$ ): A variable whose outcomes depend on randomness

Probability ( $P(X)$ ): chance of event  $X$  occurring

Expectation ( $E[X]$ ): Expected outcome over  $n$  trials

Variance ( $\text{Var}[X]$ ): Expected squared deviation of R.V. from  $E[X]$  (or mean)

Standard Deviation ( $\sigma$ ): Expected deviation from the mean

Covariance ( $\text{Cov}(X, Y)$ ): The joint variability between two R.V.

Correlation ( $\text{Corr}(X, Y)$ ): A measure of the statistical relationship between two R.V.

Probability Density Function (PDF): function that shows the distribution of values of R.V.

Cumulative Distribution Function (CDF): Integral of PDF, tells us  $P(a \leq X \leq b)$

## Formulas

$X$  is a random variable with sample space  $\Omega$

$$- E[X] = \sum_{\omega \in \Omega} P(X=\omega) \cdot \omega$$

$$E[X+Y] = E[X] + E[Y]$$

$$E[aX] = a E[X]$$

$$\begin{aligned} - \text{Var}[X] &= E[(X - E[X])^2] = E[X^2 - 2XE[X] + E[X^2]] \\ &= E[X^2] - 2(E[X])^2 + (E[X])^2 \\ &= E[X^2] - (E[X])^2 \end{aligned}$$

$$\text{Var}[aX] = a^2 \text{Var}[X]$$

$$\begin{aligned} \text{Var}[X+Y] &= \text{Var}[X] + \text{Var}[Y] + 2\text{Cov}(X, Y) \\ &= \text{Var}[X] + \text{Var}[Y] \text{ (if } \text{Cov} = 0) \end{aligned}$$

$$\text{Var}[X+b] = \text{Var}[X]$$

$$- \sigma = \sqrt{\text{Var}(X)}$$

$$\begin{aligned} - \text{Cov}(X, Y) &= E[(X - E[X])(Y - E[Y])] \\ &= E[XY] - E[X]E[Y] \end{aligned}$$

$$- \text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \cdot \text{Var}(Y)}} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$$

## Concepts / Theorems

- Common PDFs: (Discrete)

- Bernoulli:  $p = P(X=1)$   $q = 1-p = P(X=0)$

- Binomial:  $P(X) = \binom{n}{x} p^x q^{n-x}$   $n = \# \text{ trials}$   
 $x = \# \text{ success}$   
 $p = P(X=1)$   
 $q = P(X=0) = 1-p$

- Uniform:  $P(X=x) = \frac{1}{n}$   $n = \# \text{ outcomes}$

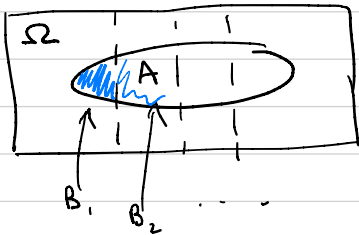
- (continuous)

- Uniform:  $X \sim U(a, b)$   $f(x) = \begin{cases} \frac{1}{b-a} & \text{if } a \leq x \leq b \\ 0 & \text{otherwise} \end{cases}$

$$P(c < x < d) = \frac{c-d}{b-a}$$

- Normal (Gaussian):  $f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$

- Law of Total Probability:



$$P(A) = P(A|B_1) \cdot P(B_1) + P(A|B_2) \cdot P(B_2) + \dots$$

$$P(A) = \sum_{i=1}^k P(A|B_i) \cdot P(B_i)$$

$P(\text{Umbrella}) = P(\text{Umbrella} | \text{Rains})$   
 $\text{Rains or it doesn't} + P(\text{Umbrella} | \text{no rain})$



- Bayes Theorem: 
$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

- Law of Large Numbers: As sample size increases the sample mean approaches the true population mean

- Central Limit Theorem: Any distribution with well defined mean and variance can be transformed into a new distribution that approaches a normal distribution.

3 coins

HT	HH	TT
$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$

1st flip is heads

$$P(HH | 1^{st} \text{ flip heads}) = \frac{P(1^{st} \text{ flip is heads} | HH) \cdot P(HH)}{P(1^{st} \text{ flip heads})}$$

$\nearrow 1$        $\frac{1}{3}$   
 $\uparrow$   
 $P(HH)$

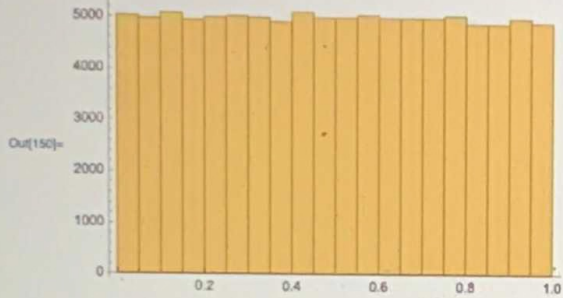
$$\begin{aligned}
 P(\text{Heads}) &= P(\text{Heads} | HT) \cdot P(HT) \\
 &\quad + P(\text{Heads} | HH) \cdot P(HH) \\
 &\quad + P(\text{Heads} | TT) \cdot P(TT) \\
 &= \frac{1}{2} \cdot \frac{1}{3} + 1 \cdot \frac{1}{3} + 0 = \frac{1}{2}
 \end{aligned}$$

$\Rightarrow \frac{\frac{1}{3}}{\frac{1}{2}} = \frac{2}{3}$



in[148]:=  $r := \text{RandomVariate}[\text{UniformDistribution}[]]$

in[150]:=  $\text{Table}[r, \{100\,000\}] // \text{Histogram}$



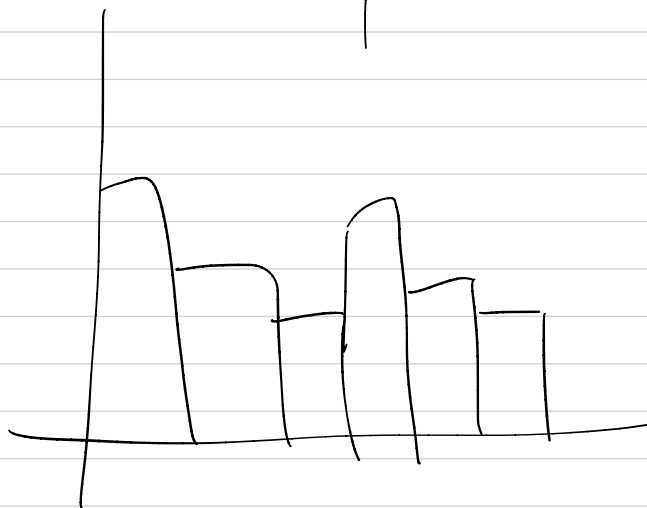
in[151]:=  $\text{Table}[r + r, \{100\,000\}] // \text{Histogram}$



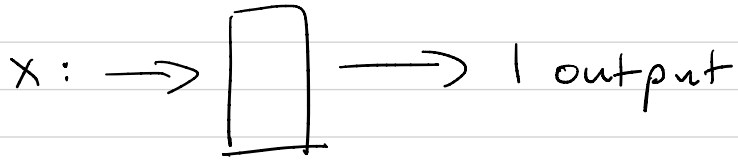
in[153]:=  $\text{Table}[r + r + r + r + r, \{100\,000\}] // \text{Histogram}$



Examples



# Neural Nets

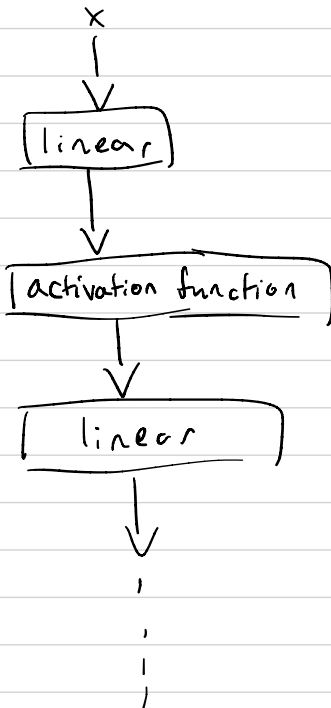
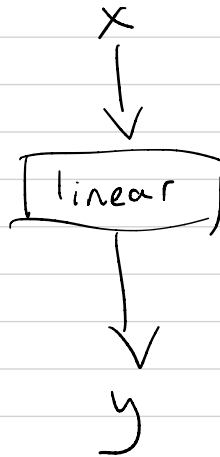

$$\frac{w}{b}$$

$$W \rightarrow \begin{bmatrix} \quad \end{bmatrix} \quad b \rightarrow \begin{bmatrix} \quad \end{bmatrix}$$

$$y = Wx + b$$

$$x = [10 \text{ values}]$$

$$y = \begin{bmatrix} 1 \\ 0 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$



$$\text{Relu}(x) = \max(0, x)$$

tanh

sigmoid

