

# MS&E 125: Intro to Applied Statistics

## Feature Engineering

Professor Udell

Management Science and Engineering  
Stanford

March 22, 2023

# Announcements

# Outline

## Linear models

To fit a linear model (= linear in parameters  $\beta$ )

- ▶ pick a transformation  $\phi : \mathcal{X} \rightarrow \mathbf{R}^p$
- ▶ predict  $y$  using a linear function of  $\phi(x)$

$$h(x) = \phi(x)^T \beta = \sum_{i=1}^p \beta_i (\phi(x))_i$$

- ▶ we want  $h(x_i) \approx y_i$  for every  $i = 1, \dots, n$

## Feature engineering

How to pick  $\phi : \mathcal{X} \rightarrow \mathbf{R}^d$ ?

- ▶ so response  $y$  will depend linearly on  $\phi(x)$
- ▶ so  $p$  is not too big

## Feature engineering

How to pick  $\phi : \mathcal{X} \rightarrow \mathbf{R}^d$ ?

- ▶ so response  $y$  will depend linearly on  $\phi(x)$
- ▶ so  $d$  is not too big

if you think this looks like a hack: you're right

## Feature engineering

examples:

- ▶ adding offset
- ▶ standardizing features
- ▶ polynomial fits
- ▶ products of features
- ▶ autoregressive models
- ▶ transforming Booleans
- ▶ transforming ordinals
- ▶ transforming nominals
- ▶ transforming images
- ▶ transforming text
- ▶ handling missing values
- ▶ concatenating data
- ▶ all of the above

<https://xkcd.com/2048/>

# Outline



## Adding offset

- ▶  $\mathcal{X} = \mathbf{R}^{d-1}$
- ▶ let  $\phi(x) = (x, 1)$
- ▶ now  $h(x) = w^T \phi(x) = w_{1:d-1}^T x + w_d$

## Fitting a polynomial

►  $\mathcal{X} = \mathbf{R}$

► let

$$\phi(x) = (1, x, x^2, x^3, \dots, x^{d-1})$$

be the vector of all monomials in  $x$  of degree  $< d$

► now  $h(x) = w^T \phi(x) = w_1 + w_2 x + w_3 x^2 + \dots + w_d x^{d-1}$

## Demo: crime

`https://colab.research.google.com/github/  
stanford-mse-125/demos/blob/main/crime.ipynb`

## Model evaluation

how should we measure how good a model is?

- ▶ (root) mean squared error (RMSE)
- ▶ mean absolute error (MAE)
- ▶ coefficient of determination ( $R^2$ )

## Mean square error

mean square error is minimized by the least squares estimator

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

equal to the sum of the residuals squared

## Root mean square error

root mean square error is the square root of the mean square error

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

(the residual standard error is similar, but normalizes by the **residual degrees of freedom**  $n - p - 1$  instead of  $n$ )

## Mean absolute error

mean absolute error is the mean of the absolute value of the residuals

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

often makes more sense than RMSE when we care about quality of the predictions

(e.g., if we will pay a linear penalty for being wrong)

## Coefficient of determination

coefficient of determination  $R^2 \in [0, 1]$  is the fraction of the variance in the data that is explained by the model

$$R^2 = 1 - \frac{\sum_i (y_i - \hat{y}_i)^2}{\sum_i (y_i - \bar{y})^2} = 1 - \frac{\text{MSE}}{\text{Var}(y)} = 1 - \frac{\text{SSR}}{\text{SST}}$$

lingo:

- ▶ SSR is the **sum of squares of the residuals**
- ▶ SST is the **total sum of squares**

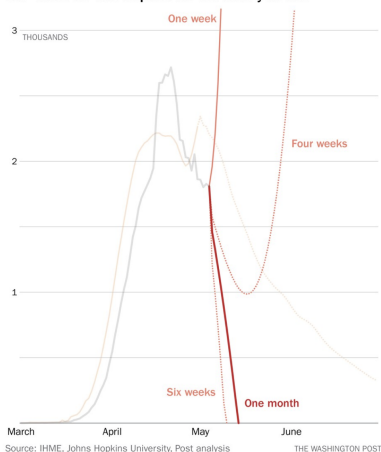
for a model with an intercept,  $R^2$  is the square correlation between the predicted and true values of  $y$

$$R^2 = [\rho(y, \hat{y})]^2$$



# IMHE and the cubic fit

The 'cubic fit' can depend on the data you use



<https://www.washingtonpost.com/politics/2020/05/05/white-houses-self-serving-approach-estimating-deadliness->

## Fitting a multivariate polynomial

- ▶  $\mathcal{X} = \mathbf{R}^2$
- ▶ pick a maximum degree  $k$
- ▶ let

$$\phi(x) = (1, x_1, x_2, x_1^2, x_1x_2, x_2^2, x_1^3, x_1^2x_2, x_1x_2^2, x_2^3, \dots, x_2^k)$$

be the vector of all monomials in  $x_1$  and  $x_2$  of degree  $\leq k$

- ▶ now  $h(x) = w^T \phi(x)$  can fit any polynomial of degree  $\leq k$  in  $\mathcal{X}$

## Fitting a multivariate polynomial

- ▶  $\mathcal{X} = \mathbf{R}^2$
- ▶ pick a maximum degree  $k$
- ▶ let

$$\phi(x) = (1, x_1, x_2, x_1^2, x_1x_2, x_2^2, x_1^3, x_1^2x_2, x_1x_2^2, x_2^3, \dots, x_2^k)$$

be the vector of all monomials in  $x_1$  and  $x_2$  of degree  $\leq k$

- ▶ now  $h(x) = w^T \phi(x)$  can fit any polynomial of degree  $\leq k$  in  $\mathcal{X}$

and similarly for  $\mathcal{X} = \mathbf{R}^d \dots$

# Outline

## Notation: boolean indicator function

define

$$\mathbb{1}(\text{statement}) = \begin{cases} 1 & \text{statement is true} \\ 0 & \text{statement is false} \end{cases}$$

examples:

- ▶  $\mathbb{1}(1 < 0) = 0$
- ▶  $\mathbb{1}(17 = 17) = 1$

## Boolean variables

- ▶  $\mathcal{X} = \{\text{true}, \text{false}\}$
- ▶ let  $\phi(x) = \mathbb{1}(x)$

## Boolean expressions

- ▶  $\mathcal{X} = \{\text{true}, \text{false}\}^2 = \{(\text{true}, \text{true}), (\text{true}, \text{false}), (\text{false}, \text{true}), (\text{false}, \text{false})\}.$
- ▶ let  $\phi(x) = [\mathbb{1}(x_1), \mathbb{1}(x_2), \mathbb{1}(x_1 \text{ and } x_2), \mathbb{1}(x_1 \text{ or } x_2)]$
- ▶ equivalent: polynomials in  $[\mathbb{1}(x_1), \mathbb{1}(x_2)]$  span the same space
- ▶ encodes logical expressions!

## Nominal values: one-hot encoding

- ▶ nominal data: e.g.,  $\mathcal{X} = \{\text{apple}, \text{orange}, \text{banana}\}$
- ▶ let

$$\phi(x) = [\mathbb{1}(x = \text{apple}), \mathbb{1}(x = \text{orange}), \mathbb{1}(x = \text{banana})]$$

- ▶ called **one-hot encoding**: only one element is non-zero



## Nominal values: one-hot encoding

- ▶ nominal data: e.g.,  $\mathcal{X} = \{\text{apple}, \text{orange}, \text{banana}\}$
- ▶ let

$$\phi(x) = [\mathbb{1}(x = \text{apple}), \mathbb{1}(x = \text{orange}), \mathbb{1}(x = \text{banana})]$$

- ▶ called **one-hot encoding**: only one element is non-zero

**extension:** sets

## Demo: crime

`https://colab.research.google.com/github/  
stanford-mse-125/demos/blob/main/crime.ipynb`

## Nominal values: the long tail

- ▶ **problem:** too many nominal categories
- ▶ **solution:**

## Nominal values: the long tail

- ▶ **problem:** too many nominal categories
- ▶ **solution:**
  - ▶ cluster the categories by some known ontology (eg, “squamous cell carcinoma” → “cancer”)

## Nominal values: the long tail

- ▶ **problem:** too many nominal categories
- ▶ **solution:**
  - ▶ cluster the categories by some known ontology (eg, “squamous cell carcinoma” → “cancer”)
  - ▶ lump the least common categories into a single category: “Other”

## Nominal values: the long tail

- ▶ **problem:** too many nominal categories
- ▶ **solution:**
  - ▶ cluster the categories by some known ontology (eg, “squamous cell carcinoma” → “cancer”)
  - ▶ lump the least common categories into a single category: “Other”
  - ▶ feature hashing

## Nominal values: the long tail

- ▶ **problem:** too many nominal categories
- ▶ **solution:**
  - ▶ cluster the categories by some known ontology (eg, “squamous cell carcinoma” → “cancer”)
  - ▶ lump the least common categories into a single category: “Other”
  - ▶ feature hashing
  - ▶ ... be creative!

## Nominal values: look up features!

why not use other information known about each item?

- ▶  $\mathcal{X} = \{\text{apple, orange, banana}\}$ 
  - ▶ price, calories, weight, ...
- ▶  $\mathcal{X} = \text{zip code}$ 
  - ▶ average income, temperature in July, walk score, % residential, ...
- ▶ ...

database lingo: **join** tables on nominal value



## Ordinal values: real encoding

- ▶ ordinal data: e.g.,  
 $\mathcal{X} = \{\text{Stage I}, \text{Stage II}, \text{Stage III}, \text{Stage IV}\}$

- ▶ let

$$\phi(x) = \begin{cases} 1, & x = \text{Stage I} \\ 2, & x = \text{Stage II} \\ 3, & x = \text{Stage III} \\ 4, & x = \text{Stage IV} \end{cases}$$

- ▶ default encoding

## Ordinal values: real encoding

- ▶  $\mathcal{X} = \{\text{Stage I, Stage II, Stage III, Stage IV}\}$
- ▶  $\mathcal{Y} = \mathbf{R}$ , number of years lived after diagnosis
- ▶ use real encoding  $\phi$  to transform ordinal data
- ▶ fit linear model with offset to predict  $y$  as  $w\phi(x) + b$

Suppose model predicts a person diagnosed with Stage II cancer will survive 2 more years, and a person diagnosed with Stage I cancer will survive 4 more years.

## Ordinal values: real encoding

- ▶  $\mathcal{X} = \{\text{Stage I, Stage II, Stage III, Stage IV}\}$
- ▶  $\mathcal{Y} = \mathbf{R}$ , number of years lived after diagnosis
- ▶ use real encoding  $\phi$  to transform ordinal data
- ▶ fit linear model with offset to predict  $y$  as  $w\phi(x) + b$

Suppose model predicts a person diagnosed with Stage II cancer will survive 2 more years, and a person diagnosed with Stage I cancer will survive 4 more years.

**Q:** What is  $w$ ?  $b$ ?

## Ordinal values: real encoding

- ▶  $\mathcal{X} = \{\text{Stage I, Stage II, Stage III, Stage IV}\}$
- ▶  $\mathcal{Y} = \mathbf{R}$ , number of years lived after diagnosis
- ▶ use real encoding  $\phi$  to transform ordinal data
- ▶ fit linear model with offset to predict  $y$  as  $w\phi(x) + b$

Suppose model predicts a person diagnosed with Stage II cancer will survive 2 more years, and a person diagnosed with Stage I cancer will survive 4 more years.

**Q:** What is  $w$ ?  $b$ ?

- A.  $b = 6, w = -2$
- B.  $b = 2, w = 0$
- C.  $b = 6, w = 2$
- D.  $b = 0, w = -2$

## Ordinal values: real encoding

- ▶  $\mathcal{X} = \{\text{Stage I, Stage II, Stage III, Stage IV}\}$
- ▶  $\mathcal{Y} = \mathbf{R}$ , number of years lived after diagnosis
- ▶ use real encoding  $\phi$  to transform ordinal data
- ▶ fit linear model with offset to predict  $y$  as  $w\phi(x) + b$

Suppose model predicts a person diagnosed with Stage II cancer will survive 2 more years, and a person diagnosed with Stage I cancer will survive 4 more years:  $b = 6$ ,  $w = -2$ .

## Ordinal values: real encoding

- ▶  $\mathcal{X} = \{\text{Stage I, Stage II, Stage III, Stage IV}\}$
- ▶  $\mathcal{Y} = \mathbf{R}$ , number of years lived after diagnosis
- ▶ use real encoding  $\phi$  to transform ordinal data
- ▶ fit linear model with offset to predict  $y$  as  $w\phi(x) + b$

Suppose model predicts a person diagnosed with Stage II cancer will survive 2 more years, and a person diagnosed with Stage I cancer will survive 4 more years:  $b = 6$ ,  $w = -2$ .

**Q:** How long does the model predict a person with Stage IV cancer will survive?

## Ordinal values: real encoding

- ▶  $\mathcal{X} = \{\text{Stage I, Stage II, Stage III, Stage IV}\}$
- ▶  $\mathcal{Y} = \mathbf{R}$ , number of years lived after diagnosis
- ▶ use real encoding  $\phi$  to transform ordinal data
- ▶ fit linear model with offset to predict  $y$  as  $w\phi(x) + b$

Suppose model predicts a person diagnosed with Stage II cancer will survive 2 more years, and a person diagnosed with Stage I cancer will survive 4 more years:  $b = 6$ ,  $w = -2$ .

**Q:** How long does the model predict a person with Stage IV cancer will survive?

- A. 6 years
- B. 2 years
- C. 0 years
- D. -2 years

## Ordinal values: boolean encoding

- ▶ ordinal data: e.g.,  
 $\mathcal{X} = \{\text{Stage I}, \text{Stage II}, \text{Stage III}, \text{Stage IV}\}$
- ▶ let

$$\phi(x) = [\mathbb{1}(x \geq \text{Stage II}), \mathbb{1}(x \geq \text{Stage III}), \mathbb{1}(x \geq \text{Stage IV})]$$



## Ordinal values: boolean encoding

- ▶  $\mathcal{X} = \{\text{Stage I, Stage II, Stage III, Stage IV}\}$
- ▶  $\mathcal{Y} = \mathbf{R}$ , number of years lived after diagnosis
- ▶ define transformation  $\phi : \mathcal{X} \rightarrow \mathbf{R}$  as

$$\phi(x) = [\mathbb{1}(x \geq \text{Stage II}), \mathbb{1}(x \geq \text{Stage III}), \mathbb{1}(x \geq \text{Stage IV})]$$

- ▶ fit linear model with offset to predict  $y$  as  $w^\top \phi(x) + b$

Suppose model predicts a person diagnosed with Stage II cancer will survive 2 more years, and a person diagnosed with Stage I cancer will survive 4 more years.

## Ordinal values: boolean encoding

- ▶  $\mathcal{X} = \{\text{Stage I, Stage II, Stage III, Stage IV}\}$
- ▶  $\mathcal{Y} = \mathbf{R}$ , number of years lived after diagnosis
- ▶ define transformation  $\phi : \mathcal{X} \rightarrow \mathbf{R}$  as

$$\phi(x) = [\mathbb{1}(x \geq \text{Stage II}), \mathbb{1}(x \geq \text{Stage III}), \mathbb{1}(x \geq \text{Stage IV})]$$

- ▶ fit linear model with offset to predict  $y$  as  $w^\top \phi(x) + b$

Suppose model predicts a person diagnosed with Stage II cancer will survive 2 more years, and a person diagnosed with Stage I cancer will survive 4 more years.

**Q:** What is  $w$ ?  $b$ ?

## Ordinal values: boolean encoding

- ▶  $\mathcal{X} = \{\text{Stage I, Stage II, Stage III, Stage IV}\}$
- ▶  $\mathcal{Y} = \mathbf{R}$ , number of years lived after diagnosis
- ▶ define transformation  $\phi : \mathcal{X} \rightarrow \mathbf{R}$  as

$$\phi(x) = [\mathbb{1}(x \geq \text{Stage II}), \mathbb{1}(x \geq \text{Stage III}), \mathbb{1}(x \geq \text{Stage IV})]$$

- ▶ fit linear model with offset to predict  $y$  as  $w^\top \phi(x) + b$

Suppose model predicts a person diagnosed with Stage II cancer will survive 2 more years, and a person diagnosed with Stage I cancer will survive 4 more years.

**Q:** What is  $w$ ?  $b$ ?

**A:**  $b = 4$ ,  $w_1 = -2$ ,  $w_2$  and  $w_3$  not determined

## Ordinal values: boolean encoding

- ▶  $\mathcal{X} = \{\text{Stage I, Stage II, Stage III, Stage IV}\}$
- ▶  $\mathcal{Y} = \mathbf{R}$ , number of years lived after diagnosis
- ▶ define transformation  $\phi : \mathcal{X} \rightarrow \mathbf{R}$  as

$$\phi(x) = [\mathbb{1}(x \geq \text{Stage II}), \mathbb{1}(x \geq \text{Stage III}), \mathbb{1}(x \geq \text{Stage IV})]$$

- ▶ fit linear model with offset to predict  $y$  as  $w^\top \phi(x) + b$

Suppose model predicts a person diagnosed with Stage II cancer will survive 2 more years, and a person diagnosed with Stage I cancer will survive 4 more years.

**Q:** What is  $w$ ?  $b$ ?

**A:**  $b = 4$ ,  $w_1 = -2$ ,  $w_2$  and  $w_3$  not determined

**Q:** How long does the model predict a person with Stage IV cancer will survive?

## Ordinal values: boolean encoding

- ▶  $\mathcal{X} = \{\text{Stage I, Stage II, Stage III, Stage IV}\}$
- ▶  $\mathcal{Y} = \mathbf{R}$ , number of years lived after diagnosis
- ▶ define transformation  $\phi : \mathcal{X} \rightarrow \mathbf{R}$  as

$$\phi(x) = [\mathbb{1}(x \geq \text{Stage II}), \mathbb{1}(x \geq \text{Stage III}), \mathbb{1}(x \geq \text{Stage IV})]$$

- ▶ fit linear model with offset to predict  $y$  as  $w^\top \phi(x) + b$

Suppose model predicts a person diagnosed with Stage II cancer will survive 2 more years, and a person diagnosed with Stage I cancer will survive 4 more years.

**Q:** What is  $w$ ?  $b$ ?

**A:**  $b = 4$ ,  $w_1 = -2$ ,  $w_2$  and  $w_3$  not determined

**Q:** How long does the model predict a person with Stage IV cancer will survive?

**A:** can't say without more information

# Outline

## Missing values

handling missing values:

- ▶ remove rows/columns with missing entries

## Missing values

handling missing values:

- ▶ remove rows/columns with missing entries
- ▶ (for time series) back-fill with most recent observed value



## Missing values

handling missing values:

- ▶ remove rows/columns with missing entries
- ▶ (for time series) back-fill with most recent observed value
- ▶ impute with mean, median, or mode

## Missing values

handling missing values:

- ▶ remove rows/columns with missing entries
- ▶ (for time series) back-fill with most recent observed value
- ▶ impute with mean, median, or mode
- ▶ fancier imputation methods (covered later in this class):  
matrix completion, copula models, deep learning, ...

## Missing values

handling missing values:

- ▶ remove rows/columns with missing entries
- ▶ (for time series) back-fill with most recent observed value
- ▶ impute with mean, median, or mode
- ▶ fancier imputation methods (covered later in this class):  
matrix completion, copula models, deep learning, ...
- ▶ add new feature: Boolean indicator  $\mathbb{1}(\text{data is missing})$ 
  - ▶ can detect if missingness is informative
  - ▶ can complement imputation method
  - ▶ can use different indicators for different kinds of missingness  
(refused, missing, illegible response, ...)

## Poll

In an ambulance dataset (data taken by instruments on board an ambulance), we want to predict if the patient died. The variable “heart rate” is sometimes missing. Is missingness

- A. informative?
- B. uninformative?

## Poll

In a weather dataset, the batteries in the instruments occasionally run out before the experimenter can replace them, leaving missing data for eg temperature, humidity, or barometric pressure. Is missingness

- A. informative?
- B. uninformative?

## Talk to your neighbor

Can you think of a dataset in which missing values would be

- ▶ informative?
- ▶ uninformative?

# Outline

## Nonlinear transformations

sometimes data is easy to predict with a simple but **nonlinear** relation, e.g.

$$\log(y) = x^T \beta$$

can transform  $x$  or (even more important)  $y$



## Nonlinear transformations

sometimes data is easy to predict with a simple but **nonlinear** relation, e.g.

$$\log(y) = x^T \beta$$

can transform  $x$  or (even more important)  $y$

hints that your data might benefit from a nonlinear transform:

- ▶  $y$  is positive and heavy-tailed? try  $y \leftarrow \log(y)$
- ▶ residuals  $r = y - x_i^T \beta$  are skewed (not normal)
  - ▶ check with quantile-quantile plot

## Nonlinear transformations

sometimes data is easy to predict with a simple but **nonlinear** relation, e.g.

$$\log(y) = x^T \beta$$

can transform  $x$  or (even more important)  $y$

hints that your data might benefit from a nonlinear transform:

- ▶  $y$  is positive and heavy-tailed? try  $y \leftarrow \log(y)$
- ▶ residuals  $r = y - x_i^T \beta$  are skewed (not normal)
  - ▶ check with quantile-quantile plot

useful nonlinear transforms:

- ▶ log, exp, quantile, ...

more systematic ways to handle nonlinearities:  
copula models, deep learning

## Log transform

**Q:** what happens if  $x$  increases by 1 in the model

$$\log(y) = \beta_0 + \beta_1 x,$$

## Log transform

**Q:** what happens if  $x$  increases by 1 in the model

$$\log(y) = \beta_0 + \beta_1 x,$$

**A:**  $\log(y)$  increases by  $\beta_1$ , so  $y$  increases by  $\exp(\beta_1)$

$$\log(y) = \beta_0 + \beta_1 x \implies y = \exp(\beta_0 + \beta_1 x)$$

$$\log(y') = \beta_0 + \beta_1(x + 1) \implies y' = \exp(\beta_0 + \beta_1(x + 1)) = \exp(\beta_0 + \beta_1 x + \beta_1) = \exp(\beta_0 + \beta_1 x) \exp(\beta_1) = y \exp(\beta_1)$$

## A convenient approximation

- ▶ for small  $x$ ,  $\exp(x) \approx 1 + x$ ,
- ▶ e.g.,  $\exp(0.01) \approx 1.01$
- ▶ if  $x$  increases by 1%, then  $y$  increases by factor of  $\exp(\beta_1/100)$
- ▶ so if  $x$  increases by 1%, then  $y$  increases by factor of  $\approx \beta_1/100 = \beta_1\%$

## Log transformations of covariates

if we instead log transform  $x$ ,  $\hat{y}$  increases by  $\beta_1/100$  for each 1% increase in  $x$ .

- ▶ e.g., if  $\beta_1 = 3$ ,  $\hat{y}$  increases by  $3/100=0.03$  units for every 1% increase in  $x$ .

if we instead log transform both  $x$  and  $y$ ,  $\hat{y}$  increases by  $\beta_1\%$  for each 1% increase in  $x$ .

- ▶ e.g., if  $\beta_1 = 3$ ,  $\hat{y}$  increases by 3% for every 1% increase in  $x$ .

log transformation results in **multiplicative** increases (rather than **additive**)

# Outline

## Location

can be given as

- ▶ latitude, longitude
- ▶ zip code
- ▶ neighborhood, county, state, country

can be transformed between these!



## Location

can be given as

- ▶ latitude, longitude
- ▶ zip code
- ▶ neighborhood, county, state, country

can be transformed between these!

which makes sense for your problem?

- ▶ does nearness matter?
- ▶ are there sharp boundaries?
- ▶ are other properties of the location (eg, mean house price or crime rate) more important?

# Outline

## Text

$\mathcal{X}$  = sentences, documents, tweets, ...

- ▶ **bag of words** model (one-hot encoding):
  - ▶ pick set of words  $\{w_1, \dots, w_d\}$
  - ▶  $\phi(x) = [\mathbb{1}(x \text{ contains } w_1), \dots, \mathbb{1}(x \text{ contains } w_d)]$
  - ▶ ignores order of words in sentence

## Text

$\mathcal{X}$  = sentences, documents, tweets, ...

- ▶ **bag of words** model (one-hot encoding):

- ▶ pick set of words  $\{w_1, \dots, w_d\}$
- ▶  $\phi(x) = [\mathbb{1}(x \text{ contains } w_1), \dots, \mathbb{1}(x \text{ contains } w_d)]$
- ▶ ignores order of words in sentence

- ▶ **pre-trained neural networks:**

- ▶ sentiment analysis: <https://medium.com/@b.terryjack/nlp-pre-trained-sentiment-analysis-1eb52a9d742c>
- ▶ Universal Sentence Encoder (USE) embedding:  
[https://colab.research.google.com/github/tensorflow/hub/blob/master/examples/colab/semantic\\_similarity\\_with\\_tf\\_hub\\_universal\\_encoder.ipynb](https://colab.research.google.com/github/tensorflow/hub/blob/master/examples/colab/semantic_similarity_with_tf_hub_universal_encoder.ipynb)
- ▶ lots of others: <https://modelzoo.co/>

## Neural networks: whirlwind primer

$$\text{NN}(x) = \sigma(W_1 \sigma(W_2 \dots \sigma(W_\ell x)))$$

- ▶  $\sigma$  is a nonlinearity applied elementwise to a vector, e.g.
  - ▶ ReLU:  $\sigma(x) = \max(x, 0)$
  - ▶ sigmoid:  $\sigma(x) = \log(1 + \exp(x))$
- ▶ each  $W$  is a matrix
- ▶ trained on very large datasets, e.g., Wikipedia, YouTube

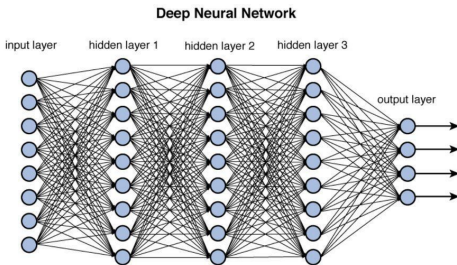


Figure 12.2 Deep network architecture with multiple layers.

# Why not use deep learning?

## Common carbon footprint benchmarks

in lbs of CO2 equivalent

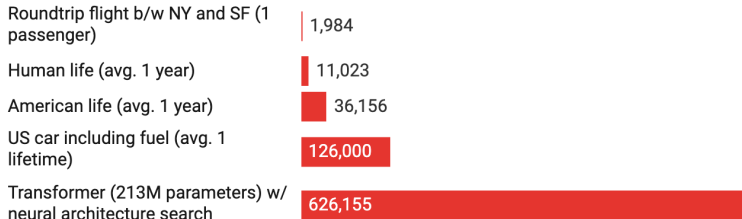


Chart: MIT Technology Review • Source: Strubell et al. • Created with Datawrapper

towards a solution: <https://arxiv.org/abs/1907.10597>

## Review

- ▶ linear models are linear in the **parameters**  $\beta$
- ▶ can fit many different models by picking feature mapping  $\phi : \mathcal{X} \rightarrow \mathbf{R}^d$