

CS-GY 6923: Lecture 2

Multiple Linear Regression + Feature Transformations + Model Selection

NYU Tandon School of Engineering, Prof. Christopher Musco

- First lab assignment `lab_housing_partial.ipynb` due **tonight, by midnight.**
- First written assignment due **next Thursday, by midnight.**
 - 10% extra credit if you use LaTeX or Markdown to typeset your assignment.

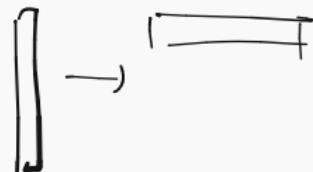
The problem set is challenging. I expect working through the problems to be one of the major ways you master material for the course. Please try to get started ASAP so that you can take advantage of office hours next week if needed.

Now it the time to review your linear algebra!

Notation:

- Let \underline{X} be an $n \times d$ matrix. Written $\underline{X} \in \mathbb{R}^{n \times d}$.
- \underline{x}_i is the i^{th} row of the matrix.
- $\underline{x}^{(j)}$ is the j^{th} column.
- x_{ij} is the i, j entry.
- For a vector \underline{y} , y_i is the i^{th} entry.
- X^T is the matrix transpose.
- y^T is a vector transpose.

$d \times n$



LINEAR ALGEBRA REVIEW

$$y \in \mathbb{R}^n$$

Things to remember:

- Matrix multiplication. If I multiply $X \in \mathbb{R}^{d \times d}$ by $Y \in \mathbb{R}^{d \times k}$ get $XY = Z \in \mathbb{R}^{n \times k}$.
- Inner product/dot product $\langle y, z \rangle = \sum_{i=1}^n y_i z_i$.
- $\langle y, z \rangle = \underline{y^T z} = z^T y$.
- Euclidean norm: $\|y\|_2 = \sqrt{y^T y}$.
- $(XY)^T = Y^T X^T$.

$$\sum_{i=1}^n y_i^2 = \|y\|_2^2 \quad \langle y, y \rangle$$

$$y^T y$$

$$\begin{bmatrix} y_1 & \dots & y_n \end{bmatrix} \xrightarrow{\text{1} \times n} \begin{bmatrix} 2, \\ \vdots \\ 2_n \end{bmatrix} \xrightarrow{n \times 1}$$

LINEAR ALGEBRA REVIEW

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

$$X\mathbb{I} = X$$

$$\mathbb{I}X = X$$

Things to remember:

- Identity matrix is denoted as \mathbb{I} . $\in \mathbb{R}^{n \times n}$
- “Most” square matrices have an inverse: i.e. if $Z \in \mathbb{R}^{n \times n}$, there is a matrix Z^{-1} such that $Z^{-1}Z = ZZ^{-1} = \mathbb{I}$.
- Let $D = \text{diag}(d)$ be a diagonal matrix containing the entries in $d = [d_1, \dots, d_n]$
- XD scales the columns of X . DX scales the rows.



LINEAR ALGEBRA REVIEW

You also need to be comfortable working with matrices in `numpy`. Go through the `demo_numpy_matrices.ipynb` slowly.

REMINDER: SUPERVISED REGRESSION

Training Dataset:

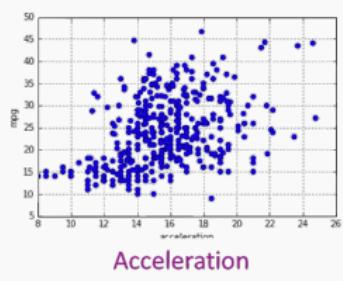
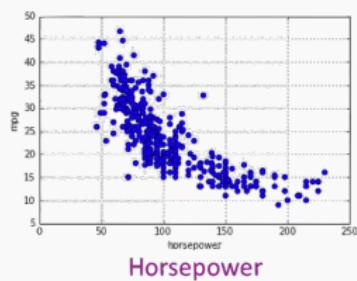
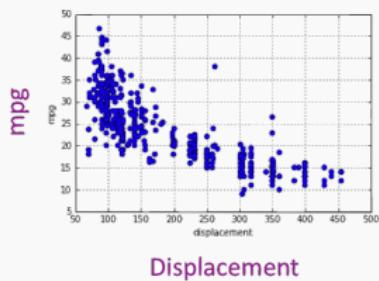
- Given input pairs $(x_1, y_1), \dots, (x_n, y_n)$.
- Each x_i is an input data vector (the predictor).
- Each y_i is a continuous output variable (the target).

Objective:

- Have the computer automatically find some function $f(x)$ such that $f(x_i)$ is close to y_i for the input data.

EXAMPLE FROM LAST CLASS

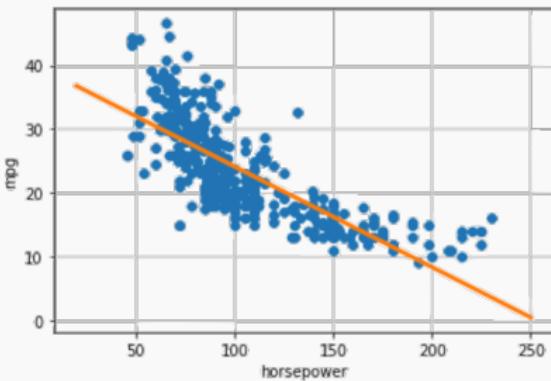
Predict miles per gallon of a vehicle given information about its engine/make/age/etc.



EXAMPLE FROM LAST CLASS

Dataset:

- $x_1, \dots, x_n \in \mathbb{R}$ (horsepowers of n cars – this is the predictor/independent variable)
- $y_1, \dots, y_n \in \mathbb{R}$ (MPG – this is the response/dependent variable)



SUPERVISED LEARNING DEFINITIONS

What are the three components needed to setup a supervised learning problem?

- **Model** $f_{\theta}(x)$: Class of equations or programs which map input x to predicted output. We want $f_{\theta}(x_i) \approx y_i$ for training inputs.
- **Model Parameters** θ : Vector of numbers. These are numerical nobs which parameterize our class of models.
- **Loss Function** $L(\theta)$: Measure of how well a model fits our data. Typically some function of $f_{\theta}(x_1) - y_1, \dots, f_{\theta}(x_n) - y_n$

Empirical Risk Minimization: Choose parameters θ^* which minimize the Loss Function:

$$\theta^* = \arg \min_{\theta} L(\theta)$$

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

Simple Linear Regression

- Model: $f_{\beta_0, \beta_1}(x) = \underline{\beta_0} + \underline{\beta_1} \cdot \underline{x}$
- Model Parameters: $\underline{\beta_0}, \underline{\beta_1}$
- Loss Function: $L(\beta_0, \beta_1) = \underbrace{\sum_{i=1}^n (y_i - f_{\beta_0, \beta_1}(x_i))^2}_{}$

Goal: Choose β_0, β_1 to minimize
 $L(\beta_0, \beta_1) = \sum_{i=1}^n |y_i - \beta_0 - \beta_1 x_i|^2.$

MULTIPLE LINEAR REGRESSION

Predict target y using multiple features, simultaneously.

Motivating example: Predict diabetes progression in patients after 1 year based on health metrics. (Measured via numerical score.)

Features: Age, sex, body mass index, average blood pressure, six blood serum measurements (e.g. cholesterol, lipid levels, iron, etc.)

Demo in `demo_diabetes.ipynb`.

LIBRARIES FOR THIS DEMO

Introducing Scikit Learn.

The screenshot displays the official scikit-learn website. At the top, there's a navigation bar with links for 'Install', 'User Guide', 'API', 'Examples', and 'More'. Below the header, the main title 'scikit-learn' is prominently displayed, followed by the subtitle 'Machine Learning in Python'. A navigation bar below the title includes 'Getting Started', 'What's New in 0.22.1', and 'GitHub'.

The page is organized into several sections:

- Classification**: Describes identifying which category an object belongs to. Applications include spam detection and image recognition. Algorithms mentioned are SVM, nearest neighbors, random forests, and more. An 'Examples' section shows a grid of plots for digit recognition.
- Regression**: Describes predicting a continuous-valued attribute associated with an object. Applications include drug response and stock prices. Algorithms mentioned are SVR, nearest neighbors, random forests, and more. An 'Examples' section shows a plot of house price prediction using a decision tree.
- Clustering**: Describes automatic grouping of similar objects into sets. Applications include customer segmentation and grouping experiment outcomes. Algorithms mentioned are k-Means, spectral clustering, mean-shift, and more. An 'Examples' section shows a scatter plot with colored regions representing different clusters.
- Dimensionality reduction**: Describes reducing the number of random variables to consider. Applications include visualization and increased efficiency. Algorithms mentioned are k-Means, feature selection, non-negative matrix factorization, and more. An 'Examples' section shows a 3D scatter plot of the Iris dataset.
- Model selection**: Describes comparing, validating, and choosing parameters and models. Applications include improved accuracy via parameter tuning. Algorithms mentioned are grid search, cross-validation, metrics, and more. An 'Examples' section shows a plot of model performance curves.
- Preprocessing**: Describes feature extraction and normalization. Applications include transforming input data such as text for use with machine learning algorithms. Algorithms mentioned are preprocessing, feature extraction, and more. An 'Examples' section shows a grid of plots for image processing tasks.



Pros:

- One of the most popular “traditional” ML libraries.
- Many built in models for regression, classification, dimensionality reduction, etc.
- Easy to use, works with ‘numpy’, ‘scipy’, other libraries we use.
- Great for rapid prototyping, testing models.

Cons:

- Everything is very “black-box”: difficult to debug, understand why models aren’t working, speed up code, etc.

Modules used:

- `datasets` module contains a number of pre-loaded datasets. Saves time over downloading and importing with `pandas`.
- `linear_model` can be used to solve Multiple Linear Regression. A bit overkill for this simple model, but gives you an idea of `sklearn`'s general structure.

THE DATA MATRIX

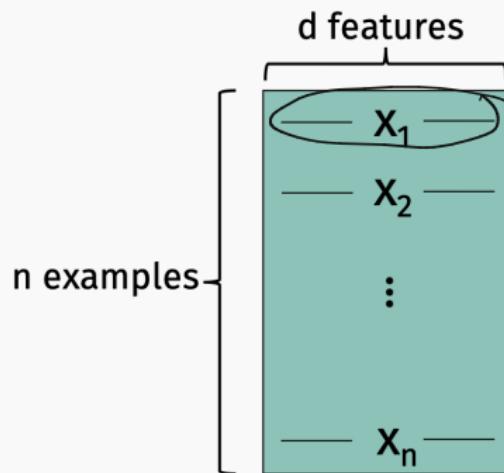
Target variable:

$n \times 1$

- Scalars y_1, \dots, y_n for n data examples (a.k.a. samples).

Predictor variables:

- d dimensional vectors x_1, \dots, x_n for n data examples and d features



THE DATA MATRIX

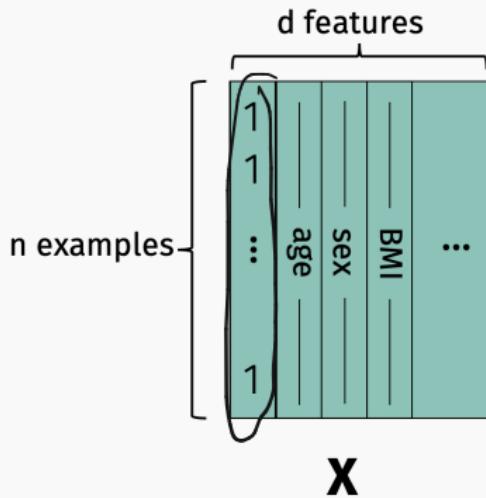
Target variable:

- Scalars y_1, \dots, y_n for n data examples (a.k.a. samples).

Predictor variables:

$$\underbrace{\beta_0}_{\text{Intercept}} + \underbrace{\beta_1}_{\text{Slope}} \times$$

- d dimensional vectors x_1, \dots, x_n for n data examples and d features



MULTIPLE LINEAR REGRESSION

Data matrix indexing:

$$Z = X_i$$
$$X = \begin{bmatrix} | & x_{11} & x_{12} & \dots & x_{1d} \\ | & x_{21} & x_{22} & \dots & x_{2d} \\ | & \textcircled{x_{31}} & \textcircled{x_{32}} & \dots & \textcircled{x_{3d}} \\ | & \vdots & \vdots & & \vdots \\ | & x_{n1} & x_{n2} & \dots & x_{nd} \end{bmatrix}$$
$$f_{\beta_0, \dots, \beta_d}(z) = \beta_0 + \beta_1 z_1 + \beta_2 z_2 + \dots + \beta_d z_d$$

Multiple Linear Regression Model:

$$= \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3}$$

Predict

$$\underline{y_i} \approx \cancel{\beta_0} + \cancel{\beta_1 x_{i1}} + \cancel{\beta_2 x_{i2}} + \dots + \cancel{\beta_d x_{id}}$$

The rate at which diabetes progresses depends on many factors, with each factor having a different magnitude effect.

MULTIPLE LINEAR REGRESSION

Assume first columns contains all 1's. If it doesn't append on a column of all 1's.

$$X = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1d} \\ x_{21} & x_{22} & \dots & x_{2d} \\ x_{31} & x_{32} & \dots & x_{3d} \\ \vdots & \vdots & & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nd} \end{bmatrix} = \begin{bmatrix} 1 & x_{12} & \dots & x_{1d} \\ 1 & x_{22} & \dots & x_{2d} \\ 1 & x_{32} & \dots & x_{3d} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n2} & \dots & x_{nd} \end{bmatrix}$$

Multiple Linear Regression Model:

Predict

$$y_i \approx \underbrace{\beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_d x_{id}}$$

MULTIPLE LINEAR REGRESSION

Use as much linear algebra notation as possible!

- Model:

$$f_{\beta_1, \dots, \beta_d}(\mathbf{z}) = \beta_1 z_1 + \beta_2 z_2 + \dots + \beta_d z_d \\ = \sum_{i=1}^d \beta_i z_i = \langle \mathbf{z}, \boldsymbol{\beta} \rangle = \mathbf{z}^\top \boldsymbol{\beta}$$

- Model Parameters:

$$\underline{\beta_1, \dots, \beta_d} \quad \hat{\boldsymbol{\beta}} \in \mathbb{R}^d \quad \begin{bmatrix} X \end{bmatrix} \quad \vec{y} = \begin{bmatrix} y \end{bmatrix}$$

- Loss Function:

$$\sum_{i=1}^n \underbrace{(y_i - (\beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_d x_{id}))^2}_{f(x_i)} = \|\vec{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2$$

Linear Least-Squares Regression.

- Model:

$$f_{\beta}(\mathbf{x}) = \langle \mathbf{x}, \boldsymbol{\beta} \rangle$$

- Model Parameters:

$$\boldsymbol{\beta} = [\beta_1, \beta_2, \dots, \beta_d]$$

- Loss Function:

$$\begin{aligned} L(\boldsymbol{\beta}) &= \sum_{i=1}^n |y_i - \langle \mathbf{x}_i, \boldsymbol{\beta} \rangle|^2 \\ &= \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 \end{aligned}$$

LINEAR ALGEBRAIC FORM OF LOSS FUNCTION

$$\sum_{i=1}^n \left(y_i - (\beta_0 x_{i0} + \beta_1 x_{i1} + \dots + \beta_d x_{id}) \right)^2 = \|y - X\beta\|_2^2$$

Diagram illustrating the components of the loss function:

- y : $n \times 1$ vector
- $X\beta$: $(n \times d) \beta \times 1$ vector
- $y - X\beta$: $n \times 1$ vector
- β : $d \times 1$ vector
- x_i : $n \times 1$ vector
- y_i : $n \times 1$ scalar

$$y_i = (\beta_0 x_{i0} + \beta_1 x_{i1} + \dots + \beta_d x_{id})$$

$$y - X\beta$$

$$\|y - X\beta\|_2^2$$

LOSS MINIMIZATION

Machine learning goal: minimize the loss function

$$L(\beta) : \mathbb{R}^d \rightarrow \mathbb{R} \quad L(\beta) = \|\gamma - X\beta\|_2^2$$

Find optimum by determining for which $\beta = [\beta_1, \dots, \beta_d]$ all partial derivatives are 0. I.e. when do we have:

gradient {

$$\begin{bmatrix} \frac{\partial L}{\partial \beta_1} \\ \frac{\partial L}{\partial \beta_2} \\ \vdots \\ \frac{\partial L}{\partial \beta_d} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \dots \\ 0 \end{bmatrix}$$

THE ALL IMPORTANT GRADIENT

For any function $L(\beta) : \mathbb{R}^d \rightarrow \mathbb{R}$, the gradient $\nabla L(\beta)$ is a function from $\mathbb{R}^d \rightarrow \mathbb{R}^d$ defined:

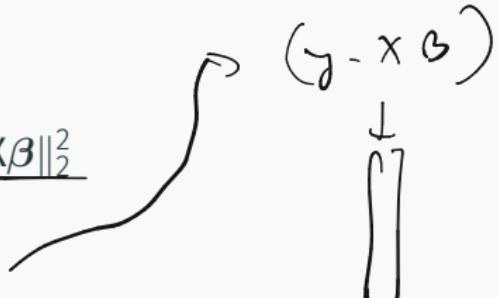
$$\nabla L(\beta) = \begin{bmatrix} \frac{\partial L}{\partial \beta_1} \\ \frac{\partial L}{\partial \beta_2} \\ \vdots \\ \frac{\partial L}{\partial \beta_d} \end{bmatrix}$$

The gradient of the loss function is a central tool in machine learning. We will use it again and again.

GRADIENT

Loss function:

$$L(\beta) = \|\underline{y} - \underline{X}\beta\|_2^2$$



Gradient:

$$\nabla L(\beta) = -2 \cdot \underline{X}^T (\underline{y} - \underline{X}\beta) = \underline{0}_{n \times 1}$$

Find optimum by determining for which $\beta = [\beta_1, \dots, \beta_d]$ the gradient is 0. i.e. when do we have:

$$\nabla L(\beta) = \begin{bmatrix} \frac{\partial L}{\partial \beta_1} \\ \frac{\partial L}{\partial \beta_2} \\ \vdots \\ \frac{\partial L}{\partial \beta_d} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \begin{array}{l} n \times d \\ (d \times n) \quad n \times 1 \\ (d \times 1) \end{array}$$

LOSS MINIMIZATION

Goal: minimize the loss function $L(\beta) = \|\mathbf{y} - \mathbf{X}\beta\|_2^2$.
 $(n \times d)(d \times k)$
 $\mathcal{O}(ndk)$

$$n \left[\begin{array}{c} d \\ \hline \end{array} \right]$$

$$\nabla L(\beta) = -2 \cdot \mathbf{X}^T (\mathbf{y} - \mathbf{X}\beta) = 0$$

$$= 2\mathbf{X}^T \mathbf{X}\beta - 2\mathbf{X}^T \mathbf{y} = 0$$

Solve for optimal β^* :

$$\cancel{\mathbf{X}^T \mathbf{X}} \beta^* = \cancel{\mathbf{X}^T \mathbf{y}}$$

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

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

$$n \left[\begin{array}{c} d \\ \hline \end{array} \right] \mathbf{X}$$

$$(1 \times n) (n \times d) \rightarrow (d \times d)^{-1} \rightarrow d \times d \rightarrow \mathcal{O}(nd^2)$$

MULTIPLE LINEAR REGRESSION SOLUTION

Need to compute $\beta^* = \arg \min_{\beta} \|y - X\beta\|_2^2 = (\underline{X^T X})^{-1} X^T y.$

- Main cost is computing $(X^T X)^{-1}$ which takes $\underline{O(nd^2)}$ time.
- Can solve slightly faster using the method $O(nd)$ `numpy.linalg.lstsq`, which is running an algorithm based on QR decomposition.
- For larger problems, can solve much faster using an iterative methods like `scipy.sparse.linalg.lsqr`.

Will learn more about iterative methods when we study
Gradient Descent.

GRADIENT WARMUP

Function:

$$\frac{\partial}{\partial z_i} \sum_{i=1}^n a_i z_i$$

$$g(z) = c \cdot z$$

$$f(z) = \underline{a^T z} \text{ for some fixed vector } a \in \mathbb{R}^d$$

Gradient:

$$g(z) \rightarrow \begin{bmatrix} \frac{\partial g}{\partial z_1} \\ \frac{\partial g}{\partial z_2} \\ \vdots \\ \frac{\partial g}{\partial z_d} \end{bmatrix} = \begin{bmatrix} a_1 \\ \vdots \\ a_n \end{bmatrix}$$

Function:

$$\underbrace{f(z) = \|z\|_2^2}_{= \sum_{i=1}^n z_i^2}$$

Gradient:

$$\nabla f(z) = \underline{2z}$$

$$\frac{\partial}{\partial z_i} \sum_{i=1}^n z_i^2 = \frac{\partial}{\partial z_i} z_i^2 = 2z_i$$

GRADIENT

Loss function:

$$X^\top \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ \vdots \\ x^{(n)} \end{bmatrix} \rightarrow X\beta - y$$

$$L(\beta) = \|y - X\beta\|_2^2$$

$$2X^\top(X\beta - y)$$

$$\frac{\partial}{\partial \beta_j} \sum_{i=1}^n (y_i - \langle x_i, \beta \rangle)^2$$

$$[x_{i1}, \dots, x_{id}]$$

$$\sum_{i=1}^n \frac{\partial}{\partial \beta_j} (y_i - \langle x_i, \beta \rangle)^2$$

$$\underline{X\beta} = \begin{bmatrix} \langle x_1, \beta \rangle \\ \langle x_n, \beta \rangle \\ \vdots \\ \langle x_i, \beta \rangle \end{bmatrix}$$

$$= \sum_{i=1}^n 2(y_i - \langle x_i, \beta \rangle) \cdot (-x_{ij})$$

$$= 2 \underbrace{\sum_{i=1}^n -y_i x_{ij}} + 2 \underbrace{\sum_{i=1}^n \langle x_i, \beta \rangle \cdot x_{ij}}$$

$$= -2 \langle y, x^{(j)} \rangle + 2 \langle X\beta, x^{(j)} \rangle$$

$$\frac{\partial}{\partial \beta_j} = 2(x^{(j)}, X\beta - y)$$

TEST YOUR INTUITION

$$\frac{\partial}{\partial \beta_j} = 2 \langle x^{(j)}, X \beta - y \rangle \quad \frac{\partial}{\partial \beta_1} = 2(x^{(1)})^T (X \beta - y)$$

$$X \quad (n \times d) \quad (d \times 1) \rightarrow (n \times 1)$$

Example from book: What is the sign of β_1 when we run a simple linear regression using the following predictors for number of sales in a particular market as a function of:

- Amount of TV advertising in that market:
- Amount of print advertising in that market:

$$\langle w, z \rangle \quad \begin{bmatrix} \frac{\partial}{\partial \beta_1} \\ \frac{\partial}{\partial \beta_2} \\ \vdots \\ \frac{\partial}{\partial \beta_d} \end{bmatrix} = 2 \begin{bmatrix} (x^{(1)})^T (X \beta - y) \\ (x^{(2)})^T (X \beta - y) \\ \vdots \\ (x^{(d)})^T (X \beta - y) \end{bmatrix} = 2x^1(X\beta - y)$$

INTERACTING VARIABLES

Positive

Positive

$$\begin{pmatrix} \beta_0 & \beta_1 \\ 0 & \beta_2 \end{pmatrix}$$

$$\begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$$

$$\begin{pmatrix} \beta_0 + \beta_1 x_1 \\ \beta_0 + \beta_2 x_2 \\ \vdots \\ \beta_0 + \beta_n x_n \end{pmatrix}$$

What is the sign of the corresponding β 's when we run a multiple linear regression using the following predictors together:

- Amount of TV advertising in that market: Positive β_1 ,
- Amount of print advertising in that market: Negative, close to zero β_2

Can you explain this? Try to think of your own example of a regression problem where this phenomenon might show up.

$$\begin{pmatrix} x^{(1)} & \dots & x^{(d)} \end{pmatrix}$$

$$\begin{pmatrix} \beta_0 & \dots & 0 \\ 0 & \dots & \beta_d \end{pmatrix}$$

$$\rightarrow \begin{pmatrix} \beta_0 + \beta_1 x^{(1)} & \beta_0 + \beta_2 x^{(2)} & \dots & \beta_0 + \beta_d x^{(d)} \end{pmatrix}$$

DEALING WITH CATEGORICAL VARIABLES

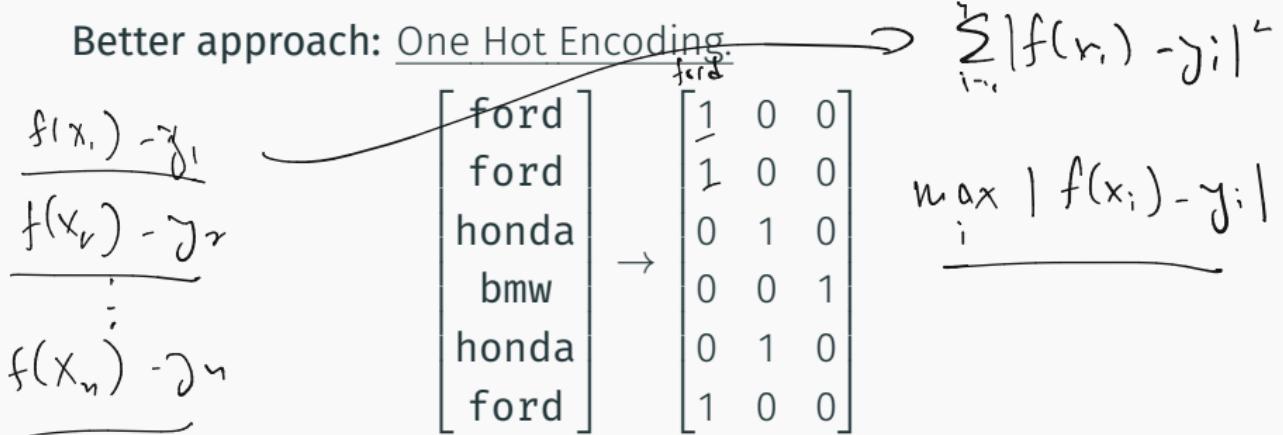
The sex variable in the diabetes problem was binary. We encoded it as 2 numbers – e.g. (0,1), (-1,1), (1,2).

Suppose we go back to the MPG prediction problem. What if we had a categorical predictor variable for car make with more than 2 options: e.g. Ford, BMW, Honda. **How would you encode as a numerical column?**

The diagram illustrates the encoding of categorical car make variables into a numerical column. On the left, a vertical vector contains categorical labels: 'ford', 'ford', 'honda', 'bmw', 'honda', and 'ford'. An arrow points from this vector to a second vertical vector on the right, which contains numerical values: 1, 1, 3, 2, 3, and 1 respectively. To the right of the second vector, handwritten labels map the numbers to car makes: 'ford' is 1, 'bmw' is 2, and 'honda' is 3.

ford	1	ford 1
ford	1	bmw 2
honda	3	honda 3
bmw	2	
honda	3	
ford	1	

ONE HOT ENCODING



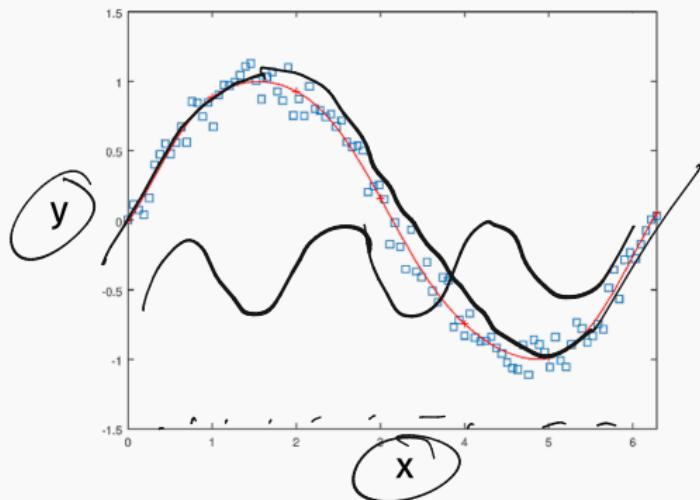
- Create a separate feature for every category, which is 1 when the variable is in that category, zero otherwise.
- Not too hard to do by hand, but you can also use library functions like sklearn.preprocessing.OneHotEncoder.

Avoids adding inadvertent linear relationships.

TRANSFORMED LINEAR MODELS

Suppose we have singular variate data examples (x, y) . How could we fit the non-linear model:

$$y \approx \underbrace{\beta_0}_{\text{ }} + \underbrace{\beta_1}_{\text{ }} x + \underbrace{\beta_2}_{\text{ }} x^2 + \underbrace{\beta_3}_{\text{ }} x^3.$$



TRANSFORMED LINEAR MODELS

Transform into a multiple linear regression problem:

$$\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \xrightarrow{x^{\star 2}} X = \begin{bmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ 1 & x_2 & x_2^2 & x_2^3 \\ 1 & x_3 & x_3^2 & x_3^3 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 & x_n^3 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} \approx \begin{bmatrix} y \end{bmatrix}$$

Each column j is generated by a different basis function $\phi_j(x)$.
Could have:

- $\phi_j(x) = \underline{x^q}$
- $\phi_j(x) = \underline{\sin(x)}$
- $\phi_j(x) = \underline{\cos(10x)}$
- $\phi_j(x) = 1/x$

TRANSFORMED LINEAR MODELS

Transformations can also be for multivariate data.

Example: Multinomial model.

- Given a dataset with target y and predictors x, z .
- For inputs $(x_1, z_1), \dots, (x_n, z_n)$ construct the data matrix:

$$\begin{bmatrix} x_1 & z_1 \\ \vdots & \vdots \\ x_n & z_n \end{bmatrix} \quad \begin{bmatrix} 1 & x_1 & x_1^2 & z_1 & z_1^2 & x_1 z_1 \\ 1 & x_2 & x_2^2 & z_2 & z_2^2 & x_2 z_2 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 & z_n & z_n^2 & x_n z_n \end{bmatrix} \quad x_1^2 z_1, \quad 2,^2 x_1$$

- Captures non-linear interaction between x and y .

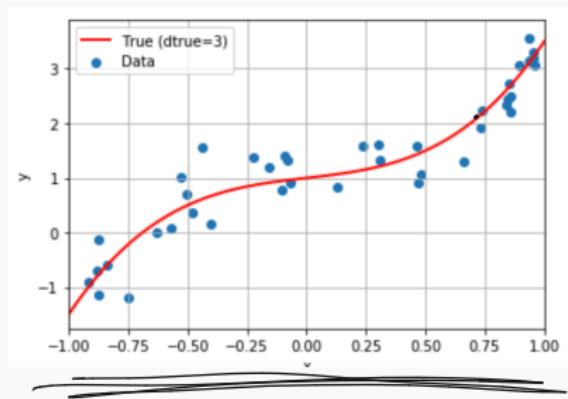
MODEL SELECTION

Remainder of lecture: Learn about model selection, test/train paradigm, and cross-validation through a simple example.

FITTING A POLYNOMIAL

Simple experiment:

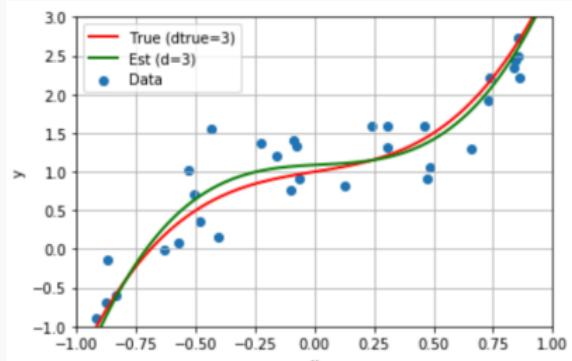
- Randomly select data points $\underline{x_1}, \dots, \underline{x_n} \in [-1, 1]$.
- Choose a degree 3 polynomial $p(x)$.
- Create some fake data: $\underline{y_i} = \underline{p(x_i)} + \eta$ where η is a random number (e.g random Gaussian).



FITTING A POLYNOMIAL

Simple experiment:

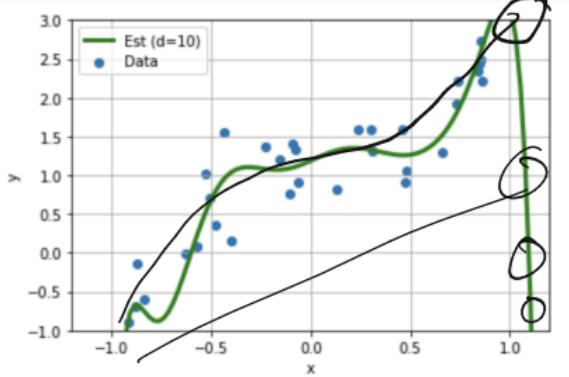
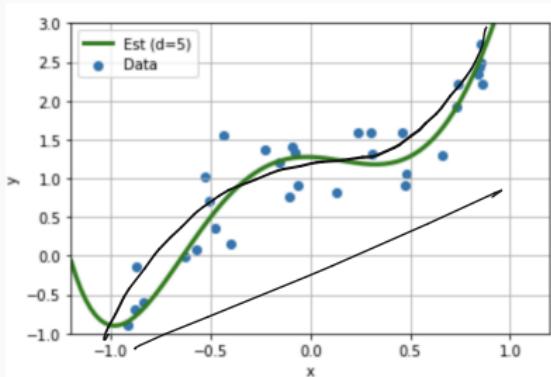
- Use multiple linear regression to fit a degree 3 polynomial.



FITTING A POLYNOMIAL

What if we fit a higher degree polynomial?

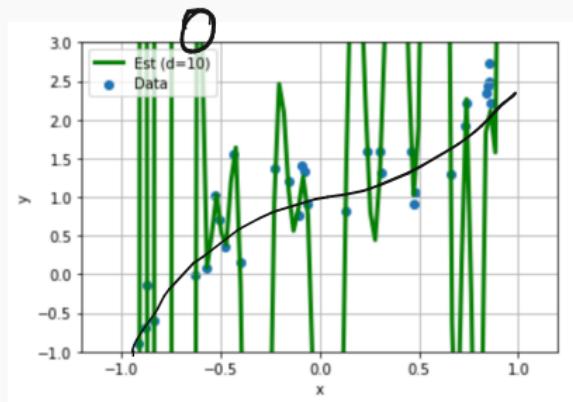
- Fit degree 5 polynomial under squared loss.
- Fit degree 10 polynomial under squared loss.



FITTING A POLYNOMIAL

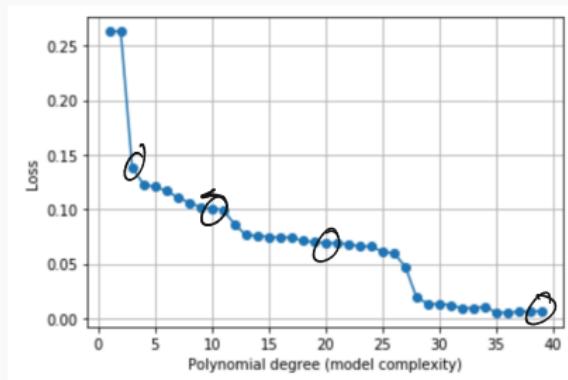
Even higher?

- Fit degree 40 polynomial under squared loss.



MODEL SELECTION

The more **complex** our model class (i.e. the higher degree we allow) the better our loss:



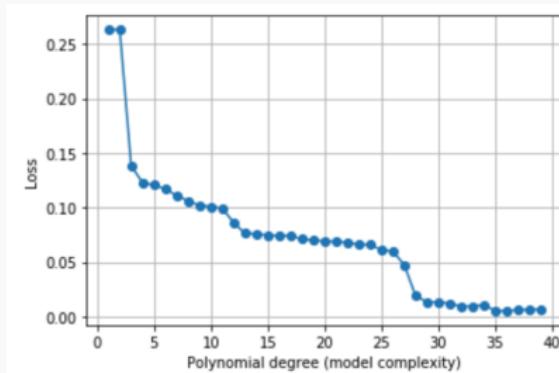
Is our model getting better and better?

Given the raw data, how do we know which model to choose?

Degree 3? Degree 5? Degree 40?

MODEL SELECTION

The more **complex** our model class the better our loss:

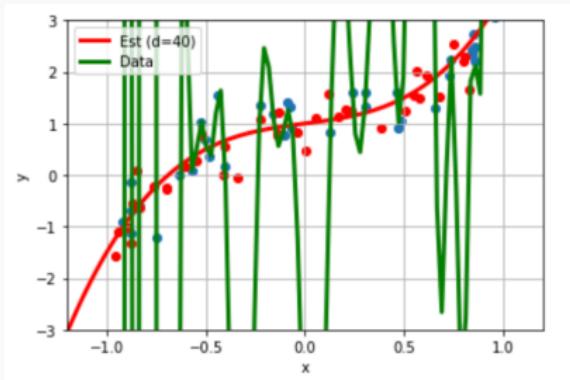


So training loss alone is not usually a good metric for model selection. Small loss does not imply generalization.

MODEL SELECTION

Problem: Loss alone is not informative for choosing model.

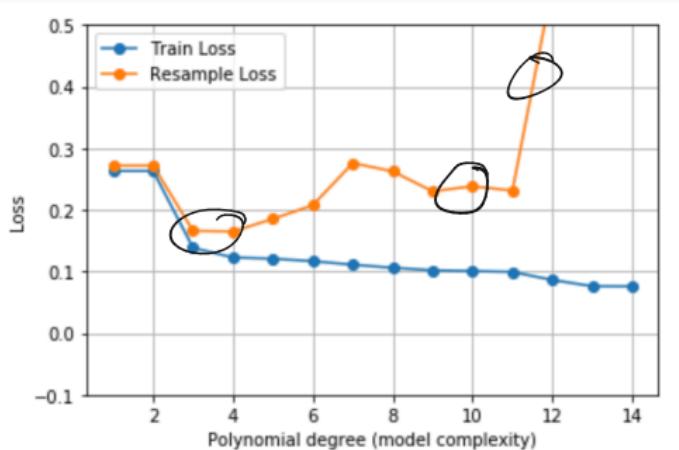
For more complex models, we get smaller loss on the training data, but don't expect to perform well on "new" data:



In other words, the model does not **generalize**.

MODEL SELECTION

Solution: Directly test model on “new data”.



- Loss continues to decrease as model complexity grows.
- Performance on new data “turns around” once our model gets too complex. Minimized around degree 4.

TRAIN-TEST PARADIGM

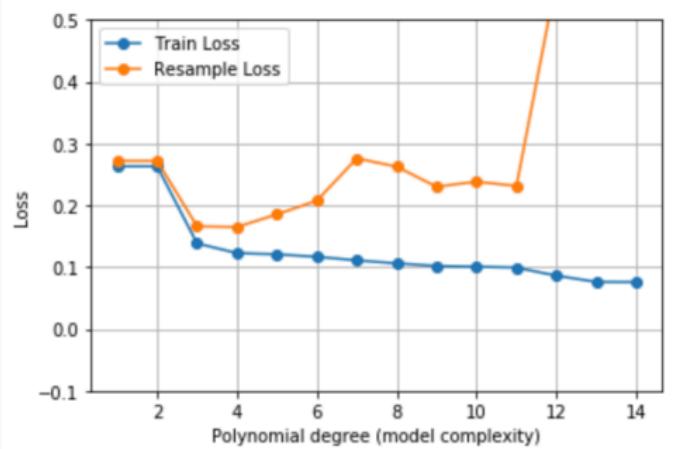
Better approach: Evaluate model on fresh test data which was not used during training.

Test/train split:

- Given data set $(\underline{X}, \underline{y})$, split into two sets $(\underline{X}_{\text{train}}, \underline{y}_{\text{train}})$ and $(\underline{X}_{\text{test}}, \underline{y}_{\text{test}})$.
- Train q models $f^{(1)}, \dots, f^{(q)}$ by finding parameters which minimize the loss on $(\underline{X}_{\text{train}}, \underline{y}_{\text{train}})$.
- Evaluate loss of each trained model on $(\underline{X}_{\text{test}}, \underline{y}_{\text{test}})$.

Sometimes you will see the term **validation set** instead of test set. Sometimes there will be both: use validation set for choosing the model, and test set for getting a final performance measure.

TRAIN-TEST PARADIGM



- **Train loss** continues to decrease as model complexity grows.
- **Test loss** “turns around” once our model gets too complex. Minimized around degree 3 – 4.

GENERALIZATION ERROR

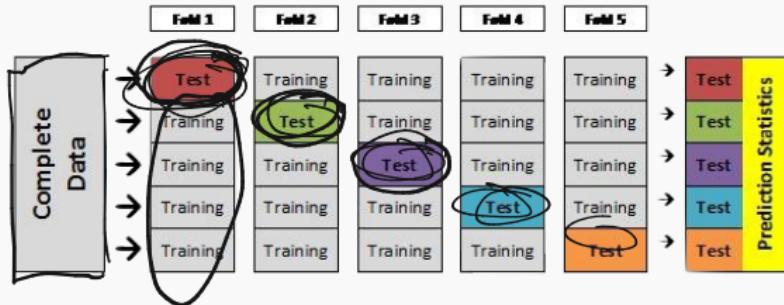
If the test loss remains low, we say that the model **generalizes**.
Test loss is often called **generalization error**.

TRAIN-TEST PARADIGM



Typical train-test split: 70-90% / 10-30%. Trade-off between between optimization of model parameters and better estimate of model performance.

K-FOLD CROSS VALIDATION



- Randomly divide data in K parts.
 - Typical choice: $K = 5$ or $K = 10$.
- Use $K - 1$ parts for training, 1 for test.
- For each model, compute test loss L_{ts} for each “fold”.
- Choose model with best average loss.
- Retrain best model on entire dataset.

K-FOLD CROSS VALIDATION

Leave-one-out cross validation: take $K = n$, where n is our total number of samples.

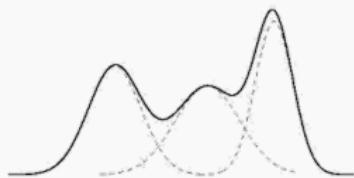
Is there any disadvantage to choosing K larger?

Intuition: Models which perform better on the test set will **generalize** better to future data.

Goal: Introduce a little bit of formalism to better understand what this means. What is “future” data?

Statistical Learning Model:

- Assume each data example is randomly drawn from some distribution $(x, y) \sim \mathcal{D}$.



This is not a simplifying assumption! The distribution could be arbitrarily complicated.

Statistical Learning Model:

- Assume each data example is randomly drawn from some distribution $(\mathbf{x}, y) \sim \mathcal{D}$.
- Define the **Risk** of a model/parameters:

$$R(f, \boldsymbol{\theta}) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [L(f(\mathbf{x}, \boldsymbol{\theta}), y)]$$

here L is our loss function (e.g. $L(z, y) = |z - y|$ or $L(z, y) = (z - y)^2$).

Goal: Find model $f \in \{f^{(1)}, \dots, f^{(q)}\}$ and parameter vector $\boldsymbol{\theta}$ to minimize the $R(f, \boldsymbol{\theta})$.

- (Population) Risk:

$$R(f, \theta) = \mathbb{E}_{(x,y) \sim \mathcal{D}} [L(f(x, \theta), y)]$$

- Empirical Risk: Draw $(x_1, y_1), \dots, (x_n, y_n) \sim \mathcal{D}$

$$R_E(f, \theta) = \frac{1}{n} \sum_{i=1}^n L(f(x_i, \theta), y_i)$$

Minimizing training loss is the same as minimizing the empirical risk on the training data.

Often called **empirical risk minimization**.

For any fixed model f and parameters θ ,

$$\mathbb{E}[R_E(f, \theta)] = R(f, \theta).$$

Only true if f and θ are chosen *without looking at the data used to compute the empirical risk*.

MODEL SELECTION

- Train q models $(f^{(1)}, \theta_1^*), \dots, (f^{(q)}, \theta_q^*)$.
- For each model, compute empirical risk $R_E(f^{(i)}, \theta_i^*)$ using test data.
- Since we assume our original dataset was drawn independently from \mathcal{D} , so is the random test subset.

No matter how our models were trained or how complex they are, $R_E(f^{(i)}, \theta_i^*)$ is an unbiased estimate of the true risk $R(f^{(i)}, \theta_i^*)$ for every i . Can use it to distinguish between models.

Slight caveat: This is typically not how machine learning or scientific discovery works in practice!

Typical workflow:

- Train a class of models.
- Test.
- Adjust class of models.
- Test.
- Adjust class of models.
- Cont...

Final model implicitly depends on test set because performance on the test set guided how we changed our model.

Popularity of ML benchmarks and competitions leads to adaptivity at a massive scale.

11 Active Competitions

	#DFDC Deepfake Detection Challenge Identify videos with facial or voice manipulations <small>Featured · Code Competition · 2 months to go · video data, online video</small>	\$1,000,000 1,595 teams
	Google QUEST Q&A Labeling Improving automated understanding of complex question answer content <small>Featured · Code Competition · 19 hours to go · text data, nlp</small>	\$25,000 1,559 teams
	Real or Not? NLP with Disaster Tweets Predict which Tweets are about real disasters and which ones are not <small>Getting Started · Ongoing · text data, binary classification</small>	\$10,000 2,657 teams
	Bengali.AI Handwritten Grapheme Classification Classify the components of handwritten Bengali <small>Research · Code Competition · a month to go · multiclass classification, image data</small>	\$10,000 1,194 teams

Kaggle (various competitions)



14,197,122 images, 21841 synsets indexed

[Explore](#) [Download](#) [Challenges](#) [Publications](#) [Updates](#) [About](#)

Not logged in. [Login](#) | [Signup](#)

Imagenet (image classification and categorization)

Is adaptivity a problem? Does it lead to over-fitting? How much? How can we prevent it? All current research.

REPORT

The reusable holdout: Preserving validity in adaptive data analysis

Cynthia Dwork^{1,*}, Vitaly Feldman^{2,*}, Moritz Hardt^{3,*}, Toniann Pitassi^{4,*}, Omer Reingold^{5,*}, Aaron Roth^{6,*}

* See all authors and affiliations

Science 07 Aug 2015:
Vol. 349, Issue 6248, pp. 636-638
DOI: 10.1126/science.aaa9375

Do ImageNet Classifiers Generalize to ImageNet?

Benjamin Recht*
UC Berkeley

Rebecca Roelofs
UC Berkeley

Ludwig Schmidt
UC Berkeley

Vaishaal Shankar
UC Berkeley

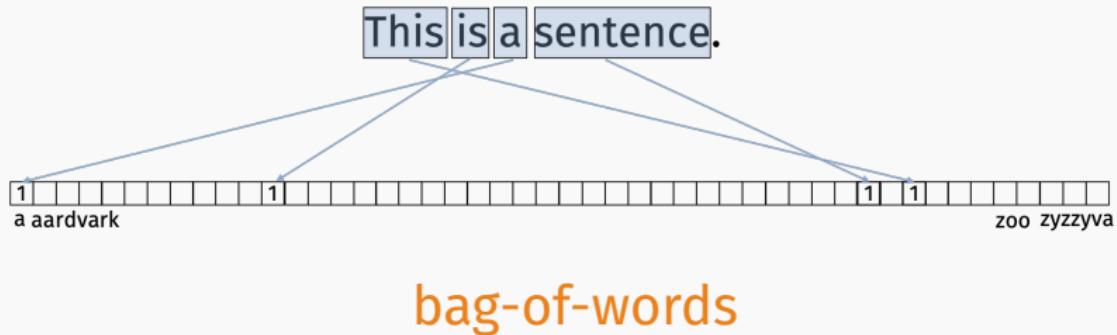
Abstract

We build new test sets for the CIFAR-10 and ImageNet datasets. Both benchmarks have been the focus of intense research for almost a decade, raising the danger of overfitting to excessively re-used test sets. By closely following the original dataset creation processes, we test to what extent current classification models generalize to new data. We evaluate a broad range of models and find accuracy drops of 3% – 15% on CIFAR-10 and 11% – 14% on ImageNet. However, accuracy gains on the original test sets translate to larger gains on the new test sets. Our results suggest that the accuracy drops are not caused by adaptivity, but by the models' inability to generalize to slightly "harder" images than those found in the original test sets.

MODEL SELECTION EXAMPLE

bag-of-words models and n-grams

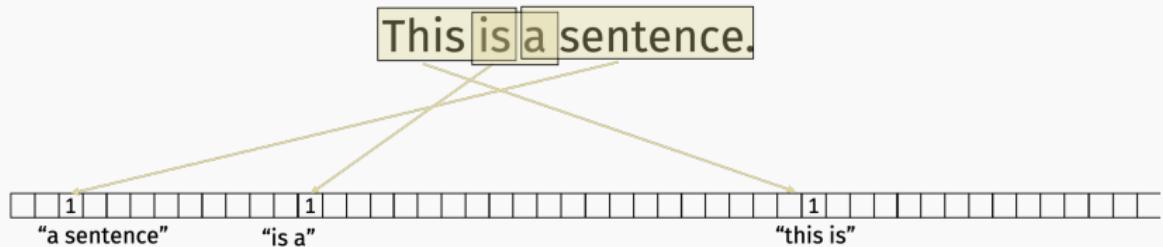
Common way to represent documents (emails, webpages, books) as numerical data. The ultimate example of 1-hot encoding.



MODEL SELECTION EXAMPLE

bag-of-words models and n-grams

Common way to represent documents (emails, webpages, books) as numerical data. The ultimate example of 1-hot encoding.

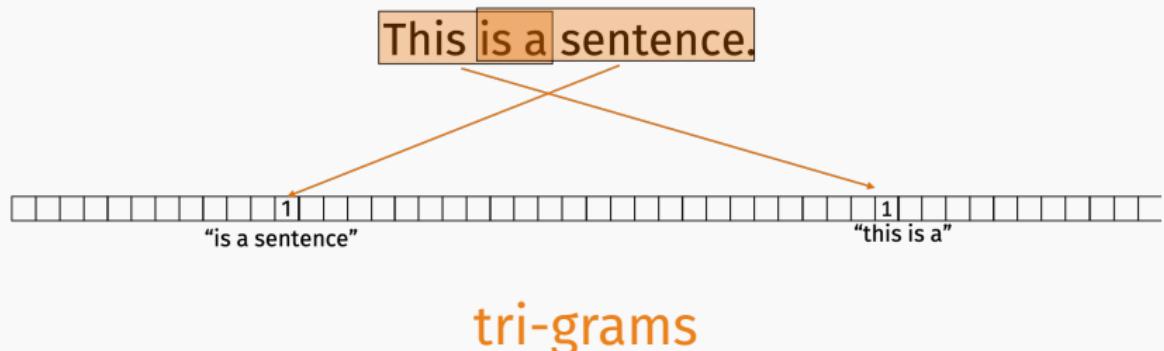


bi-grams

MODEL SELECTION EXAMPLE

bag-of-words models and n-grams

Common way to represent documents (emails, webpages, books) as numerical data. The ultimate example of 1-hot encoding.



MODEL SELECTION EXAMPLE

Models of increasing order:

- Model $f_{\theta_1}^{(1)}$: spam filter that looks at **single words**.
- Model $f_{\theta_2}^{(2)}$: spam filter that looks at **bi-grams**.
- Model $f_{\theta_3}^{(3)}$: spam filter that looks at **tri-grams**.
- ...

“interest”

“low interest”

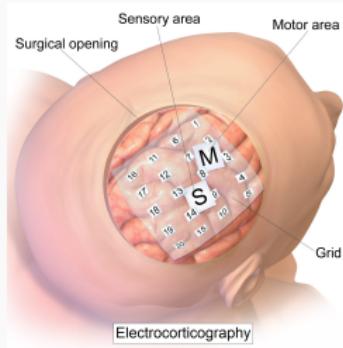
“low interest loan”

Increased length of **n-gram** means more expressive power.

MODEL SELECTION EXAMPLE

Electrocorticography ECoG (upcoming lab):

- Implant grid of electrodes on surface of the brain to measure electrical activity in different regions.



- Predict hand motion based on ECoG measurements.
- Model order:** predict movement at time t using brain signals at time $t, t - 1, \dots, t - q$ for varying values of q .