## CS-GY 6923: Lecture 2 Multiple Linear Regression + Feature Transformations + Model Selection

NYU Tandon School of Engineering, Prof. Christopher Musco

- First lab assignment lab1.ipynb due Monday, by midnight.
- First written assignment will be released this weekend.
- TA's will start office hours next week thanks for everyone who filled out the poll.

## Training Dataset:

- Given input pairs  $(\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_n, y_n)$ .
- Each  $\mathbf{x}_i$  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  $f(\mathbf{x})$  such that  $f(\mathbf{x}_i)$  is close to  $y_i$  for the input data.

**Standard approach:** Convert the supervised learning problem to a multi-variable <u>optimization problem</u>.

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



## 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  $\theta$ : 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  $\theta^*$  which minimize the Loss Function:

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

#### SIMPLE LINEAR REGRESSION

### Simple Linear Regression

- Model:  $f_{\beta_0,\beta_1}(x) = \beta_0 + \beta_1 \cdot x$
- Model Parameters:  $\beta_0, \beta_1$
- Loss Function:  $L(\beta_0, \beta_1) = \sum_{i=1}^{n} (y_i f_{\beta_0, \beta_1}(x_i))^2$
- $\mathcal{F}_{\boldsymbol{b},\boldsymbol{l}}^{\boldsymbol{a}} \mathcal{F}_{\boldsymbol{b}}^{\boldsymbol{a}}$  **Goal:** Choose  $\beta_0, \beta_1$  to minimize  $L(\beta_0, \beta_1) = \sum_{i=1}^{n} |y_i - \beta_0 - \beta_1 x_i|^2.$

Simple closed form solution:  $\beta_1 = \sigma_{xy}/\sigma_{xy}^2$ ,  $\beta_0 = \bar{y} - \beta_1 \bar{x}$ . How did we solve for this solution?

#### MULTIPLE LINEAR REGRESSION



#### MULTIPLE LINEAR REGRESSION

Model Parameters:

## Linear Least-Squares Regression.

Y: x B, X: + Bo

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

$$\underbrace{\beta}_{d_1} = [\beta_1, \beta_2, \dots, \beta_d]$$

$$\underbrace{\beta}_{d_1} = [\beta_1, \beta_2, \dots, \beta_d]$$

$$\underbrace{\beta}_{d_1} = [\beta_1, \beta_2, \dots, \beta_d]$$

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

• Loss Function:

· Model:

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

### LINEAR ALGEBRAIC FORM OF LOSS FUNCTION

$$\begin{split} \mathcal{J} - \mathcal{X} \mathcal{B} \\ \mathcal{X} \mathcal{B} &= \begin{cases} \mathcal{X}_{1}, \mathcal{B} \mathcal{I} \\ \vdots \\ \mathcal{X}_{N}, \mathcal{B} \mathcal{I} \end{cases} & \mathcal{J} = \begin{cases} \mathcal{Y}_{1} \\ \vdots \\ \mathcal{Y}_{N}, \mathcal{B} \mathcal{I} \end{cases} \\ \mathcal{Y} = \begin{cases} \mathcal{Y}_{1} - \mathcal{X}_{1}, \mathcal{B} \mathcal{I} \\ \mathcal{Y}_{N} \end{cases} \\ \mathcal{Y} = \begin{cases} \mathcal{Y}_{1} - \mathcal{X}_{1}, \mathcal{B} \mathcal{I} \\ \mathcal{Y}_{N} \end{cases} \\ \mathcal{Y} = \begin{cases} \mathcal{Y}_{1} - \mathcal{X}_{1}, \mathcal{B} \mathcal{I} \\ \mathcal{Y}_{N} \end{cases} \\ \mathcal{Y} = \begin{cases} \mathcal{Y}_{1} - \mathcal{X}_{1}, \mathcal{B} \mathcal{I} \\ \mathcal{Y}_{N} \end{cases} \\ \mathcal{Y} = \begin{cases} \mathcal{Y}_{1} - \mathcal{X}_{1}, \mathcal{B} \mathcal{I} \\ \mathcal{Y}_{N} \end{cases} \\ \mathcal{Y} = \begin{cases} \mathcal{Y}_{1} - \mathcal{X}_{1}, \mathcal{B} \mathcal{I} \\ \mathcal{Y}_{N} \end{cases} \\ \mathcal{Y} = \begin{cases} \mathcal{Y}_{1} - \mathcal{X}_{1}, \mathcal{B} \mathcal{I} \\ \mathcal{Y}_{N} \end{cases} \\ \mathcal{Y} = \begin{cases} \mathcal{Y}_{1} - \mathcal{X}_{1}, \mathcal{B} \\ \mathcal{Y}_{N} \end{cases} \\ \mathcal{Y} = \begin{cases} \mathcal{Y}_{1} - \mathcal{X}_{1}, \mathcal{B} \\ \mathcal{Y}_{N} \end{cases} \\ \mathcal{Y} = \begin{cases} \mathcal{Y}_{1} - \mathcal{X}_{1}, \mathcal{B} \\ \mathcal{Y}_{N} \end{cases} \\ \mathcal{Y} = \begin{cases} \mathcal{Y}_{1} - \mathcal{X}_{1}, \mathcal{B} \\ \mathcal{Y}_{N} \end{cases} \\ \mathcal{Y} = \begin{cases} \mathcal{Y}_{1} - \mathcal{X}_{1}, \mathcal{B} \\ \mathcal{Y}_{N} \end{cases} \\ \mathcal{Y} = \begin{cases} \mathcal{Y}_{1} - \mathcal{X}_{1}, \mathcal{B} \\ \mathcal{Y}_{N} \end{cases} \\ \mathcal{Y} = \begin{pmatrix} \mathcal{Y}_{1} - \mathcal{Y}_{N}, \mathcal{B} \\ \mathcal{Y}_{N} \end{pmatrix} \\ \mathcal{Y} = \begin{pmatrix} \mathcal{Y}_{1} - \mathcal{Y}_{N}, \mathcal{B} \\ \mathcal{Y}_{N} \end{pmatrix} \\ \mathcal{Y} = \begin{pmatrix} \mathcal{Y}_{1} - \mathcal{Y}_{N}, \mathcal{B} \\ \mathcal{Y}_{N} \end{pmatrix} \\ \mathcal{Y} = \begin{pmatrix} \mathcal{Y}_{1} - \mathcal{Y}_{N}, \mathcal{B} \\ \mathcal{Y}_{N} \end{pmatrix} \\ \mathcal{Y} = \begin{pmatrix} \mathcal{Y}_{1} - \mathcal{Y}_{N}, \mathcal{B} \\ \mathcal{Y}_{N} \end{pmatrix} \\ \mathcal{Y} = \begin{pmatrix} \mathcal{Y}_{1} - \mathcal{Y}_{N}, \mathcal{B} \\ \mathcal{Y}_{N} \end{pmatrix} \\ \mathcal{Y} = \begin{pmatrix} \mathcal{Y}_{1} - \mathcal{Y}_{N}, \mathcal{B} \\ \mathcal{Y}_{N} \end{pmatrix} \\ \mathcal{Y} = \begin{pmatrix} \mathcal{Y}_{1} - \mathcal{Y}_{N}, \mathcal{B} \\ \mathcal{Y}_{N} \end{pmatrix} \\ \mathcal{Y} = \begin{pmatrix} \mathcal{Y}_{1} - \mathcal{Y}_{N}, \mathcal{B} \\ \mathcal{Y}_{N} \end{pmatrix} \\ \mathcal{Y} = \begin{pmatrix} \mathcal{Y}_{1} - \mathcal{Y}_{N}, \mathcal{B} \\ \mathcal{Y}_{N} \end{pmatrix} \\ \mathcal{Y} = \begin{pmatrix} \mathcal{Y}_{1} - \mathcal{Y}_{N}, \mathcal{B} \\ \mathcal{Y}_{N} \end{pmatrix} \\ \mathcal{Y} = \begin{pmatrix} \mathcal{Y}_{1} - \mathcal{Y}_{N}, \mathcal{B} \\ \mathcal{Y}_{N} \end{pmatrix} \\ \mathcal{Y} = \begin{pmatrix} \mathcal{Y}_{1} - \mathcal{Y}_{N}, \mathcal{B} \\ \mathcal{Y}_{N} \end{pmatrix} \\ \mathcal{Y} = \begin{pmatrix} \mathcal{Y}_{1} - \mathcal{Y}_{N}, \mathcal{B} \\ \mathcal{Y}_{N} \end{pmatrix} \\ \mathcal{Y} = \begin{pmatrix} \mathcal{Y}_{1} - \mathcal{Y}_{N}, \mathcal{Y}_{N} \end{pmatrix} \\ \mathcal{Y} = \begin{pmatrix} \mathcal{Y}_{1} - \mathcal{Y}_{N}, \mathcal{Y}_{N} \end{pmatrix} \\ \mathcal{Y} = \begin{pmatrix} \mathcal{Y}_{1} - \mathcal{Y}_{N}, \mathcal{Y}_{N} \end{pmatrix} \\ \mathcal{Y} = \begin{pmatrix} \mathcal{Y}_{1} - \mathcal{Y}_{N}, \mathcal{Y}_{N} \end{pmatrix} \\ \mathcal{Y} = \begin{pmatrix} \mathcal{Y}_{1} - \mathcal{Y}_{N} \end{pmatrix} \\ \mathcal{Y} = \begin{pmatrix} \mathcal{Y}_{1} - \mathcal{Y}_{N} \end{pmatrix} \\ \mathcal{Y} = \begin{pmatrix} \mathcal{Y}_{1} - \mathcal{Y}_{N} \end{pmatrix} \end{pmatrix} \\ \mathcal{Y} = \begin{pmatrix} \mathcal{Y}_{1} - \mathcal{Y}_{N} \end{pmatrix} \\ \mathcal{Y} = \begin{pmatrix} \mathcal{Y}_{1} - \mathcal{Y}_{N} \end{pmatrix} \\ \mathcal{Y} = \begin{pmatrix} \mathcal{Y}_{1} - \mathcal{Y}_{N} \end{pmatrix} \end{pmatrix} \\ \mathcal{Y} = \begin{pmatrix} \mathcal{Y}_{1} - \mathcal{Y}_{N} \end{pmatrix} \end{pmatrix} \\ \mathcal{Y} = \begin{pmatrix} \mathcal{Y}_{1} - \mathcal{Y}_{N} \end{pmatrix} \end{pmatrix}$$

# **Machine learning goal:** minimize the loss function $L(\beta) : \mathbb{R}^d \to \mathbb{R}$ .

Find possible optima by determining for which  $\beta = [\beta_1, \dots, \beta_d]$ all the **gradient** equals **0**. I.e. when do we have:

$$\nabla L(\boldsymbol{\beta}) = \begin{pmatrix} \frac{\partial L}{\partial \beta_1} \\ \frac{\partial L}{\partial \beta_2} \\ \vdots \\ \frac{\partial L}{\partial \beta_d} \end{pmatrix} = \begin{bmatrix} 0 \\ 0 \\ \cdots \\ 0 \end{bmatrix}$$

GRADIENT



Can check that this is equal to 0 if 
$$\beta = (X^T X)^{-1} X^T y$$
. There are  
no other options, so this must be the minimum.  
 $\chi^T$  -  $\chi^T \chi \mathcal{B} = O$   $\chi^T \mathcal{J} - \chi^T \chi (\chi^T \chi)^{-1} \chi^T \mathcal{J}$   
 $\chi^T \mathcal{J} - \chi^T \chi (\chi^T \chi)^{-1} \chi^T \mathcal{J}$ 

11

#### SINGLE VARIABLE WARMUP

#### GRADIENT

$$\begin{array}{c} \text{Loss function:} \qquad \mathcal{B} = \begin{bmatrix} \mathcal{B}_{1}, \dots, \mathcal{B}_{d} \end{bmatrix} \qquad \mathcal{P} \lfloor (\mathcal{B}) := \begin{bmatrix} \mathcal{B}_{1} \\ \mathcal{B}_{d} \end{bmatrix} \\ \mathcal{P} \lfloor (\mathcal{B}) := \| \mathbf{y} - \mathbf{X} \mathcal{B} \|_{2}^{2} \qquad \mathcal{P} \lfloor (\mathcal{B}) := \begin{bmatrix} \mathcal{B}_{1} \\ \mathcal{B}_{d} \end{bmatrix} \\ \mathcal{P} \lfloor (\mathcal{B}) := \| \mathbf{y} - \mathbf{X} \mathcal{B} \|_{2}^{2} \qquad \mathcal{P} \lfloor (\mathcal{B}) := \begin{bmatrix} \mathcal{B}_{1} \\ \mathcal{B}_{d} \end{bmatrix} \\ \mathcal{P} \lfloor \mathcal{B} \rfloor \\ \mathcal{P} \lfloor \mathcal{P} \lfloor \mathcal{P} \rfloor \\ \mathcal{P} \lfloor \mathcal{P} \rfloor \\ \mathcal{P} \lfloor \mathcal{P} \lfloor \mathcal{P} \rfloor \\ \mathcal{P} \lfloor \mathcal{P} \lfloor \mathcal{P} \rfloor \\ \mathcal{P} \lfloor \mathcal{P} \rfloor \\ \mathcal{P} \lfloor \mathcal{P$$

Take away: simple form for the gradient means that multiple linea regression models are easy and efficient to optimize.

$$\boldsymbol{\beta}^* = \arg\min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

• Often the "go to" first regression method. Throw your data in a matrix and see what happens.

It is not always immediately clear how to do this! One of the first issue we run into is categorical data:

$$\begin{array}{c} x_{1} = [42, 4, 104, hybrid, ford] \\ x_{2} = [18, 8, 307, gas, bmw] \\ x_{2} = [31, 4, 150, gas, honda] \\ \vdots \\ & & & & \\ \end{array}$$

Binary data is easy to deal with – pick one category to be 0, one to be 1. The choice doesn't matter – it will not impact the overall loss of the model

$$x_1 = [42, 4, 104 \text{ hybrid}] \text{ford}]$$
  
 $x_2 = [18, 8, 307] \text{gas} \text{[bmw]}$   
 $x_2 = [31, 4, 150] \text{gas}, \text{honda}]$ 

 $\begin{aligned} \mathbf{x}_1 &= [42, 4, 104, \mathbf{1}, \mathbf{ford}] \\ \mathbf{x}_2 &= [18, 8, 307, \mathbf{0}, \mathbf{bmw}] \\ \mathbf{x}_2 &= [31, 4, 150, \mathbf{0}, \mathbf{honda}] \end{aligned}$ 

What about 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?



### ONE HOT ENCODING

### Better approach: One Hot Encoding.

[ford]	$\rightarrow$	1	0	0
ford		1	0	0
honda		0	1	0
bmw		0	0	1
honda		0	1	0
ford		1	0	0

- 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

Instead of fitting the model mpg  $\approx \beta_0 + \beta_1 \cdot$  horsepower, fit the model mpg  $\approx \beta_0 + \beta_1 \cdot 1/$ horsepower.



How would you know to make such a transformation?

**Better approach:** Choose a more flexible non-linear model class.

#### TRANSFORMED LINEAR MODELS



**Claim:** This can be done using an algorithm for multivariate regression! No need to compute another gradient or write good to optimize  $\beta_0, \ldots, \beta_3$ .

#### TRANSFORMED LINEAR MODELS

$$\underbrace{X}_{l} = \begin{bmatrix}
1 & x_{1} & x_{1}^{2} & x_{1}^{3} \\
1 & x_{2} & x_{2}^{2} & x_{2}^{3} \\
1 & x_{3} & x_{3}^{2} & x_{3}^{3} \\
\vdots & \vdots & \vdots \\
1 & x_{n} & x_{n}^{2} & x_{n}^{3}
\end{bmatrix} \begin{pmatrix}
\mathcal{B}_{\bullet} \\
\mathcal{B}_{\bullet} \\
\mathcal{B}_{\bullet} \\
\mathcal{B}_{\bullet}
\end{pmatrix} = \begin{pmatrix}
\mathcal{B}_{\bullet} + \mathcal{B}_{\bullet} \times \mathcal{B}_{\bullet} \times \mathcal{B}_{\bullet} \times \mathcal{B}_{\bullet} \\
\mathcal{B}_{\bullet} \times \mathcal{B}_{\bullet}$$

What is the output of  $X\beta$ ?

More generally, have each column *j* is generated by a different basis function  $\phi_i(x)$ . Could have:



When might you want to include sins and cosines?

When might you want to include sins and cosines? Time series data:



There is usually not much harm in including irrelevant variable transformation.

### Transformations can also be for multivariate data.

Example: Multinomial model.

- Given a dataset with <u>target y</u> and <u>predictors x, z</u>.
- For inputs  $(\underline{x}_1, \underline{z}_1), \ldots, (x_n, z_n)$  construct the data matrix:

$$\begin{bmatrix} 1 & x_1 & x_1^2 & z_1 & z_1^2 & x_1 z_1 \\ 1 & x_2 & x_2^2 & z_2 & z_2^2 & x_2 z_2 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 & z_n & z_n^2 & x_n z_n \end{bmatrix} \quad X_1^{\flat} Z_1$$

• Captures non-linear interaction between x and z.

We use these a lot in my work to fit models for physical phenomenon over low-dimensional surfaces:



Between at 12:20pm.

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

I have a Python demo working through this example.

$$a^{t}b = \sum_{i=1}^{2} a_{i}, b_{i}$$
  
 $b_{i}^{T}y = \sum_{i=1}^{2} y_{i}y_{i}, = \sum_{i=1}^{2} y_{i}^{*}, = \|y\|_{2}^{2}$ 

## Simple experiment:

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



## Simple experiment:

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



## What if we fit a higher degree polynomial?

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



## Even higher?

• Fit degree 40 polynomial under squared loss.



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



#### MODEL SELECTION

Consider  $X \in \mathbb{R}^{n \times d}$  and  $\overline{X} = [X, z] \in \mathbb{R}^{n \times d+1}$  with one additional column appended on. XB': X6+ 200 - XB+ Claim:  $\min_{\bar{\boldsymbol{\beta}} \in \mathbb{R}^{d+1}} \|\bar{\mathbf{X}}\bar{\boldsymbol{\beta}} - \mathbf{y}\|_2^2 \le \min_{\boldsymbol{\beta} \in \mathbb{R}^d} \|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|_2^2.$ Get opt for right hand side: 15th (our up with some B' such that B'= B\* || X b'- J || = ||xb+- J ||2 = XB\*

#### MODEL SELECTION

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



So <u>training loss</u> alone is not usually a good metric for model selection. Small loss does not imply generalization. Generalization: How well do we do on <u>new</u> data. 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.
**More reasonable approach:** Evaluate model on fresh <u>test data</u> which was not used during training.

# Test/train split:

- Given data set (X, y), split into two sets (X<sub>train</sub>, y<sub>train</sub>) and (X<sub>test</sub>, y<sub>test</sub>).
- Train q models  $f^{(1)}, \ldots, f^{(q)}$  by finding parameters which minimize the loss on  $(X_{train}, y_{train})$ .
- Evaluate loss of each trained model on  $(X_{test}, y_{test})$ .

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.

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

# **Typical train-test split:** 90-70% / 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 Lts for each "fold".
- Choose model with best average loss.
- Retrain best model on entire dataset.

### Is there any disadvantage to choosing *K* larger?

The above trend is fairly representative of what we tend to see across the board:



# Is "test error" the end goal though? Don't we care about "future" error?

(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

# Statistical Learning Model:



E.g.  $x_1, \ldots, x_d$  are Gaussian random variables with parameters  $\mu_1, \sigma_1, \ldots, \mu_d, \sigma_d$ .

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

0

## Statistical Learning Model:

$$f(x, \phi) = f_{\phi}(x)$$

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

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

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

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

$$\mathbb{E}\left(R_{E}(f,\sigma)\right) = R(f,\sigma) \qquad \mathbb{E}\left(A + B\right) = \mathbb{E}\left(A\right) + \mathbb{E}\left(A\right)$$

· (Population) Risk:

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

• Empirical Risk: Draw 
$$(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n) \sim \mathcal{D}$$
  

$$R_{\underline{F}}(f, \theta) = \underbrace{\frac{1}{n} \sum_{i=1}^{n} L(f(\mathbf{x}_i; \theta), y_i)}_{\substack{i=1 \\ j \in \mathbb{N}}} = \mathcal{P}(f, \theta)$$

$$= \underbrace{\frac{1}{n} \left[ \mathcal{F} \left( \sum_{i=1}^{n} L(f(\mathbf{x}_i; \theta), y_i) \right) - \frac{1}{n} - \sum_{i=1}^{n} \mathcal{F} \left( L(f(\mathbf{x}_i; \theta), y_i) \right) - \frac{1}{n} - \sum_{i=1}^{n} \mathcal{F} \left( f(\mathbf{x}_i; \theta), y_i \right) - \frac{1}{n} - \sum_{i=1}^{n} \mathcal{F} \left( f(\mathbf{x}_i; \theta), y_i \right) - \frac{1}{n} - \sum_{i=1}^{n} \mathcal{F} \left( f(\mathbf{x}_i; \theta), y_i \right) - \frac{1}{n} - \sum_{i=1}^{n} \mathcal{F} \left( f(\mathbf{x}_i; \theta), y_i \right) - \frac{1}{n} - \sum_{i=1}^{n} \mathcal{F} \left( f(\mathbf{x}_i; \theta), y_i \right) - \frac{1}{n} - \frac$$



Only true if f and  $\theta$  are chosen without looking at the data used to compute the empirical risk.

- Train q models  $(f^{(1)}, \boldsymbol{\theta}_1^*), \ldots, (f^{(q)}, \boldsymbol{\theta}_q^*).$
- 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 <u>unbiased estimate</u> of the true risk  $R(f^{(i)}, \theta_i^*)$  for every *i*. Can use it to distinguish between models.

## 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.



bag-of-words



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



bi-grams

## 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.



tri-grams

• . . .

# 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.

Will also be relevant in our first generative ML lab!

### 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, ..., t q for varying values of q.

Predicting time t based on a linear function of the signals at time t, t - 1, ..., t - q is not the same as fitting a line to the time series. It's much more expressive.



Predecessor of modern "recurrent neural networks".

#### MODEL SELECTION LAB TIP

## Electrocorticography ECoG lab:



First lab where computation actually matters (solving regression problems with  $\sim 40k$  examples,  $\sim 1500$  features)

Makes sense to test and debug code using a subset of the data.

#### ADAPTIVE DATA ANALYSIS

# 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.

#### ADAPTIVE DATA ANALYSIS

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

11 Active Competitions		
#DFDC	Deepfake Detection Challenge Isenthy video with tacai or voice merputations Fautured - Code Competition - 2 months to go - % video data, ontine video	\$1,000,000 1,595 teams
6	Google QUEST Q&A Labeling           Improving automated understanding of complex question answer content           Failured         - Code Competition - 18 hours to go - % lixet data, np	\$25,000 1,559 teams
	Real or Not? NLP with Disaster Tweets Predict which Tweets are about real disasters and which ones are not Getting Started - Origong - % text data, binary classification	\$10,000 2,657 teams
В	Bengali LAI Handwritten Grapheme Classification Classify the components of handwritten Bengal Research - Code Competition - a month to go - % multiclass classification, image data	\$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 I 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. Related to the problem of "p-value hacking" in science.

REPORT

# The reusable holdout: Preserving validity in adaptive data analysis

Cynthia Dwork<sup>1,\*</sup>, Vitaly Feldman<sup>2,\*</sup>, Moritz Hardt<sup>3,\*</sup>, Toniann Pitassi<sup>4,\*</sup>, Omer Reingold<sup>5,\*</sup>, Aaron Roth<sup>6,\*</sup>

+ 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<sup>\*</sup> UC Berkeley Rebecca Roelofs UC Berkeley

Ludwig Schmidt UC Berkeley Vaishaal Shankar UC Berkeley

Abstract

We build new test sets for the CIPAR-10 and ImageNet datasets. Both benchmarks have been the focus of inteme 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, 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.



Collected by Fei-Fei Li's group at Stanford in 2006ish and labeled using Amazon Mechanical Turk.



We now have neural network models that can solve these classification problems with > 95% accuracy.

### ADAPTIVE DATA ANALYSIS

### Do ImageNet Classifiers Generalized to ImageNet?



Interestingly, when comparing popular vision models on "fresh" data, while performance dropped across the board, the relative rank of model performance did not change significantly.

### REGULARIZATION

# In all the model selection examples we discussed we had full control over the complexity of the model: could range from <u>underfitting</u> to <u>overfitting</u>.

In practice, you often don't have this freedom. Even the <u>most</u> <u>basic model</u> might lead to overfitting.

**Example:** Linear regression model where  $d \ge n$ . Almost always the case e.g. when using bag-of-words features.



Can (almost) always find  $\beta$  so that  $X\beta = y$  exactly.

**Claim:** For almost all sets of *n* length *n* vectors  $\mathbf{x}^{(1)}, \ldots, \mathbf{x}^{(n)}$ , we can write any vector  $\mathbf{y}$  as a linear combination of these vectors.



### FEATURE SELECTION

### Select some subset of features to use in model:



**Filter method:** Compute some metric for each feature, and select features with highest score.

• Example: compute loss or  $R^2$  value when each feature in X is used in single variate regression.

# Any potential limitations of this approach?

### FEATURE SELECTION

Exhaustive approach: Pick best subset of q features.

Faster approach: Greedily select q features.

Stepwise Regression:

- Forward: Step 1: pick single feature that gives lowest loss.
   Step k: pick feature that when combined with previous k 1 chosen features gives lowest loss.
- **Backward:** Start with all of the features. Greedily eliminate those which have least impact on model performance.

Feature selection deserves more than two slides, but we won't go into too much more detail!

**Regularization:** Explicitly discourage overfitting by adding a regularization penalty to the loss minimization problem.

 $\min_{\boldsymbol{\theta}} \left[ L(\boldsymbol{\theta}) + Reg(\boldsymbol{\theta}) \right].$ 

**Example:** Least squares regression.  $L(\beta) = ||X\beta - y||_2^2$ .

- Ridge regression ( $\ell_2$ ):  $Reg(\beta) = \lambda \|\beta\|_2^2$
- LASSO (least absolute shrinkage and selection operator) ( $\ell_1$ ):  $Reg(\beta) = \lambda \|\beta\|_1$
- Elastic net:  $Reg(\beta) = \lambda_1 \|\beta\|_1 + \lambda_2 \|\beta\|_2^2$

### REGULARIZATION



Ridge regression:  $\min_{\beta} \|\mathbf{X}\beta - \mathbf{y}\|_{2}^{2} + \lambda \|\beta\|_{2}^{2}$ .

- As  $\lambda \to \infty$ , we expect  $\|\beta\|_2^2 \to 0$  and  $\|\mathbf{X}\beta \mathbf{y}\|_2^2 \to \|\mathbf{y}\|_2^2$ .
- Feature selection methods attempt to set many coordinates in β to 0. Ridge regularizations encourages coordinates to be small.



Ridge regression:  $\min_{\beta} \|\mathbf{X}\beta - \mathbf{y}\|_{2}^{2} + \lambda \|\beta\|_{2}^{2}$ .

• Can be viewed as shrinking the size of our model class. Relaxed version of  $\min_{\beta:||\beta||_2^2 < c} ||\mathbf{X}\beta - \mathbf{y}||_2^2$ .

**Claim:** For any  $\lambda$ , let  $\beta_{\lambda}^* = \arg \min_{\beta} \|\mathbf{X}\beta - \mathbf{y}\|_2^2 + \lambda \|\beta\|_2^2$ . Then there is some  $c(\lambda)$  such that:

$$\boldsymbol{\beta}_{\lambda}^{*} = \operatorname*{arg\,min}_{\boldsymbol{\beta}:\|\boldsymbol{\beta}\|_{2}^{2} < c(\lambda)} \|\boldsymbol{X}\boldsymbol{\beta} - \boldsymbol{y}\|_{2}^{2}.$$

Moreover, we have the for  $\lambda' > \lambda$ ,  $c(\lambda') < c(\lambda)$ .

Ridge regression:  $\min_{\beta} \|\mathbf{X}\beta - \mathbf{y}\|_{2}^{2} + \lambda \|\beta\|_{2}^{2}$ .

•  $\min_{\beta:\|\beta\|_2^2 < c} \|\mathbf{X}\beta - \mathbf{y}\|_2^2$  won't have zero error solution for all **y**, even when over-parameterized.


## How do we minimize: $L_R(\beta) = \|\mathbf{X}\beta - \mathbf{y}\|_2^2 + \lambda \|\beta\|_2^2$ ?

Lasso regularization:  $\min_{\beta} \|\mathbf{X}\beta - \mathbf{y}\|_{2}^{2} + \lambda \|\beta\|_{1}$ .

- As  $\lambda \to \infty$ , we expect  $\|\beta\|_1 \to 0$  and  $\|\mathbf{X}\beta \mathbf{y}\|_2^2 \to \|\mathbf{y}\|_2^2$ .
- Typically encourages subset of β<sub>i</sub>'s to go to zero, in contrast to ridge regularization.



## Pros:

- Simpler, more interpretable model.
- More intuitive reduction in model order.

## Cons:

- No closed form solution because  $\|\beta\|_1$  is not differentiable.
- Can be solved with iterative methods, but generally not as quickly as ridge regression.

## Notes:

- Model selection/cross validation used to choose optimal scaling  $\lambda$  on  $\lambda \|\beta\|_2^2$  or  $\lambda \|\beta\|_1$ .
- Often grid search for best parameters is performed in "log space". E.g. consider  $[\lambda_1, \ldots, \lambda_q] = 1.5^{[-4, -3, -2, -1, -0, 1, 2, 3, 4]}$ .
- Regularization methods are <u>not invariant</u> to data scaling. Typically when using regularization we mean center and scale columns to have unit variance.