

ML for SS: Workflow and ML regression

Session 1

Dr. Richard M. Crowley
rcrowley@smu.edu.sg
<http://rmc.link/>

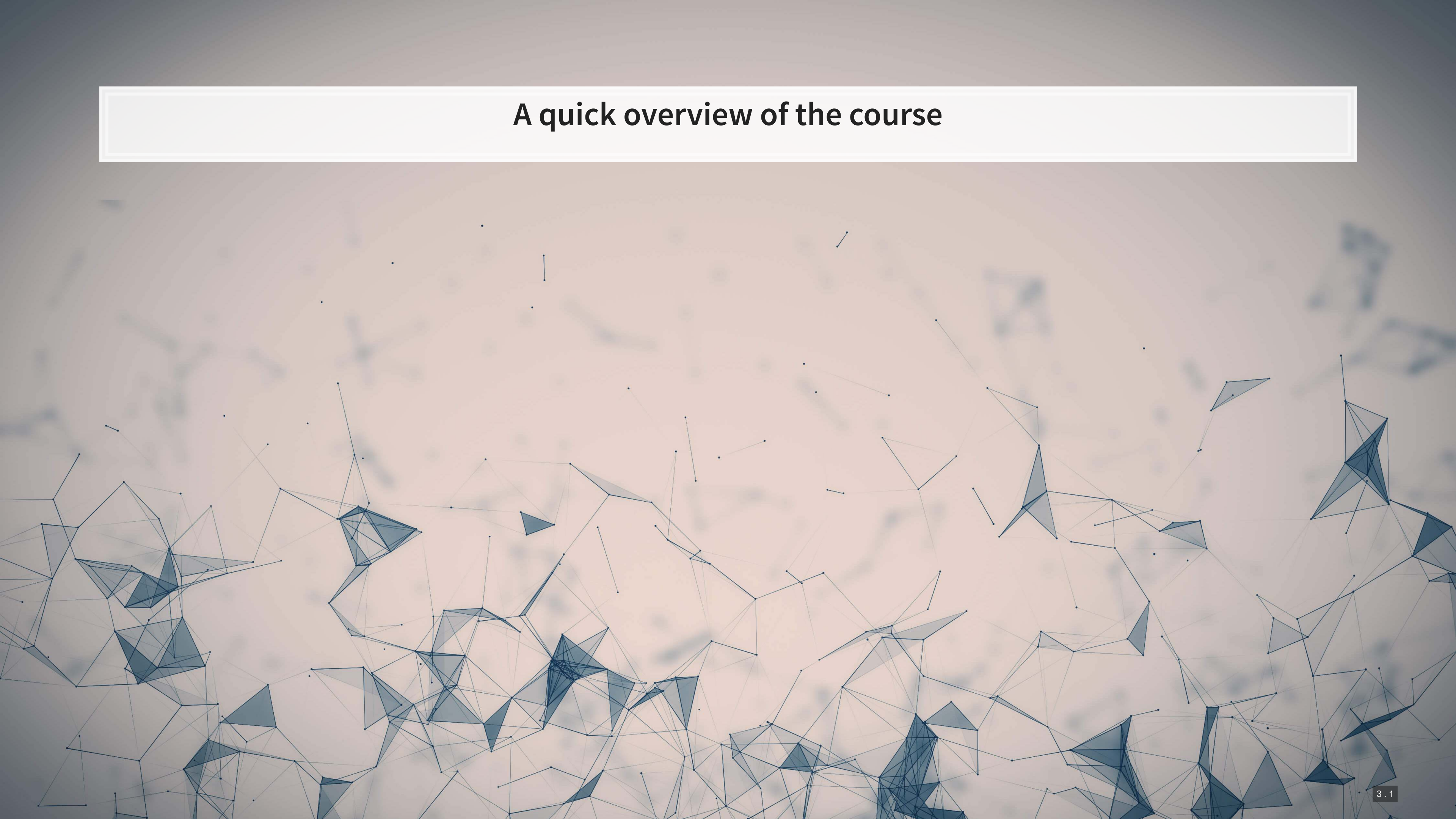
About me

Research highlights

1. An advanced model for detecting financial misreporting using the topic modeling applied to annual report text.
2. Multiple projects on Twitter showcasing:
 1. How companies are more likely to disclose both good and bad information than what is normal or expected
 2. That CSR disclosure on Twitter is not credible
 3. That executives' disclosures are as important on Twitter as their firms' disclosures
3. Newer work on
 - COVID-19 reactions worldwide
 - Sentiment and understandability in accounting text
 - Misinformation laws (e.g., POFMA)

The above all use some sort of machine learning or text analytics approach.

A quick overview of the course



A typical class session

2-3 papers to discuss

- Each paper will usually use a different method
- Almost all papers are applied
- Student led
 - 2 students per paper

Method overview

- Walk through methods' technical aspects
- Discuss how and where the method is useful
- Showcase a coded up example
 - When feasible, I will show this for both R and python
- Professor led

What we will cover

Regression and analysis with ML

1. Working with data and ML regression
2. Tree-based ML algorithms
3. Clustering algorithms

Not far removed from traditional econometrics, but more flexible

- Drop-in replacements for regression
- Non-linear and non-parametric methods
- Dimensionality reduction

What we will cover

Working with textual data

4. Text processing (NLP)
5. Linguistics
6. Embedding and topic models
7. Inferring traits from text

These methods are often useful in measuring phenomenon

- What is being discussed (content)
- People's sentiment or emotion toward something

What we will cover

Economics approaches to ML

8. Causal machine learning
9. Policy prediction
10. Bias

Useful for measuring impact or effects

- Understanding policy impacts
- Understanding processes

What we will cover

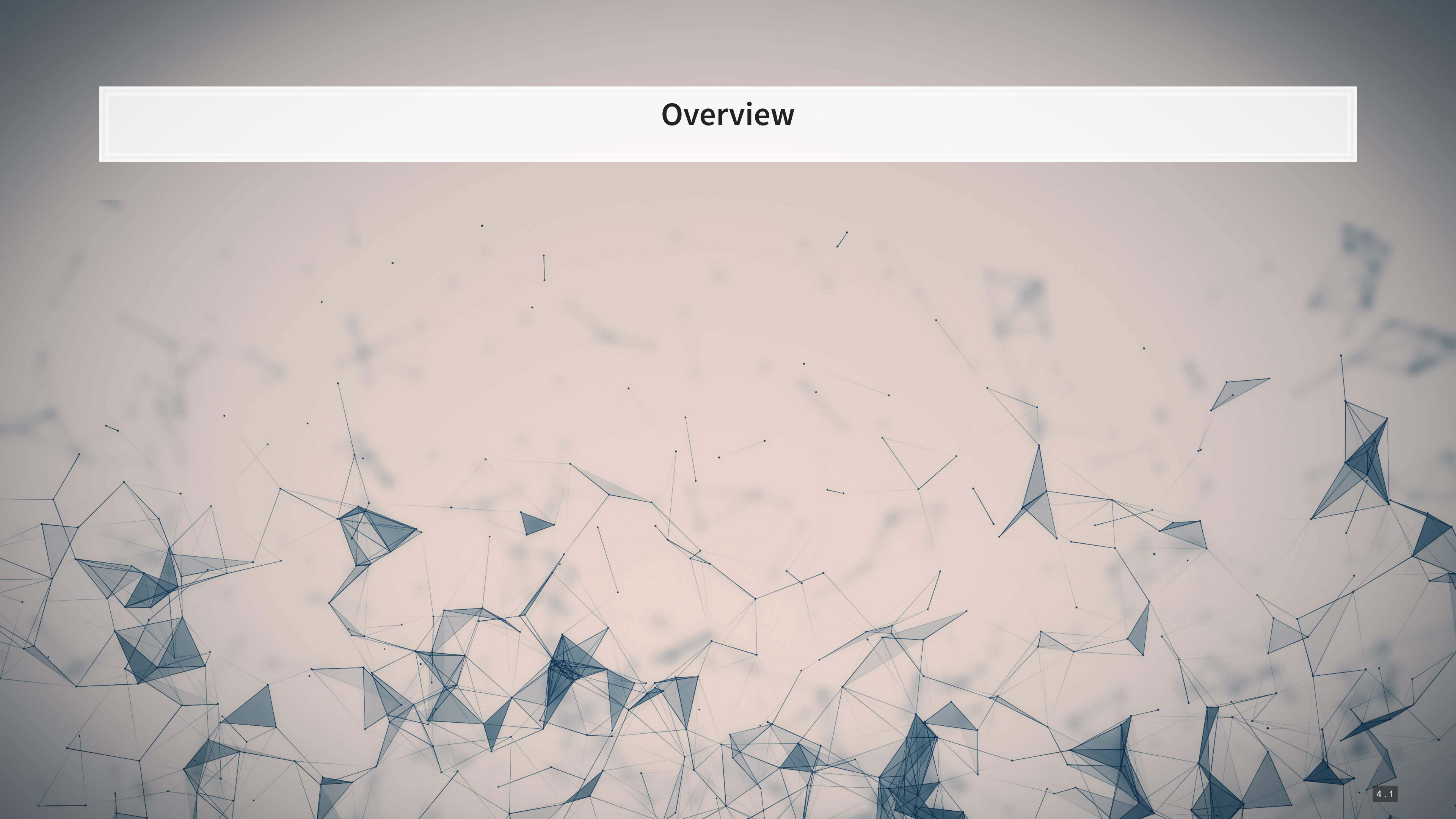
Neural networks

11. Text processing
12. Image processing

These methods can offer powerful methods for measuring phenomenon

- Better understanding message content
- Picking apart images
- Building better classifiers

Overview



Papers

Paper 1: Mullainathan and Spiess 2017 JEP

- A fairly approachable overview of ML methods in economics
- The points the paper makes are applicable broadly in any archival/empirical discipline

Paper 2: Chahuneau et al 2012

- An application of LASSO to a context most should be familiar with: restaurant menus
- Easy to motivate LASSO in this paper – more variables than observations!

Technical discussion: Implementing LASSO

1. Sample splitting
2. Cross validation
3. What are LASSO and Elastic Net
4. Implementing them

Python

- Using `sklearn`
- Can be done using built-in CV methods

R

- Using `glmnet`
- Easy to use
- Fast
- Nice CV method for penalty

Both Python and R are good for this. Stata is also pretty good with `lassopack`.

There is a fully worked out solution for each language on my website, data is on eLearn.

Main application: A linear problem

- Idea: Discussion of risks, such as as foreign currency risks, operating risks, or legal risks should provide insight on the volatility of future outcomes for the firm.
- Testing: Predicting future stock return volatility based on 10-K filing discussion

Dependent Variable

- Future stock return volatility

Independent Variables

- A set of 31 measures of what was discussed in a firm's annual report

This test mirrors Bao and Datta (2014 MS)

Secondary application: A binary problem

- Idea: Using the same data as in Application 1, can we predict instances of intentional misreporting?
- Testing: Predicting 10-K/A irregularities using finance, textual style, and topics

Dependent Variable

Intentional misreporting as stated in 10-K/A filings

Independent Variables

- 17 Financial measures
- 20 Style characteristics
- 31 10-K discussion topics

This test mirrors a subset of Brown, Crowley and Elliott (2020 JAR)

Paper 1: An overview of applied ML

Paper 2: ML for panel data

Problems of the usual approach

- For both linear and logistic regression:
 - Easy to have too many covariates
 - Which can lead to high VIFs and multicollinearity
- For logit:
 - Convergence is iffy when using sparse datasets or DVs

How can machine learning help?

1. Some methods directly address the issues of multicollinearity or having too many covariates (via model selection)
2. Some methods address sparsity well, being robust to binary DVs with sub 10% classes

What is LASSO?

- Least Absolute Shrinkage and Selection Operator
 - Least absolute: uses an error term like $|\varepsilon|$
 - Shrinkage: it will make coefficients smaller
 - Less sensitive \rightarrow less overfitting issues
 - Selection: it will completely remove some variables
 - Less variables \rightarrow less overfitting issues
- Sometimes called L_1 regularization
 - L_1 means 1 dimensional distance, i.e., $|\varepsilon|$

Great if you have way too many inputs in your model or high multicollinearity

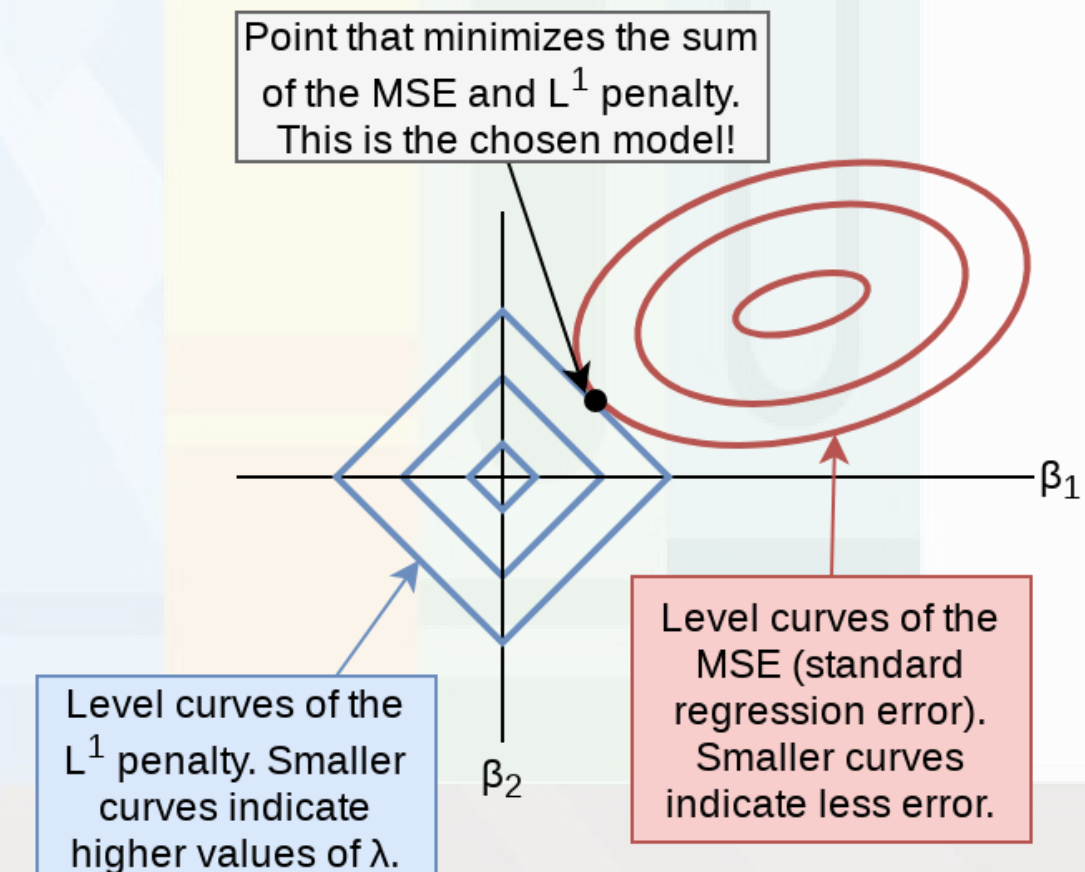
- Note that L_1 regularization is a standard approach to dealing with inflated VIFs as well!

How does it work?

$$\min_{\beta \in \mathbb{R}} \left\{ \frac{1}{N} |\varepsilon|_2^2 + \lambda |\beta|_1 \right\}$$

- Add an additional penalty term that is increasing in the absolute value of each β
 - Incentivizes lower β s, *shrinking* them
- The selection part is explainable geometrically in 2D
 - If the MSE level curves hit a corner of the diamond shaped penalty curve, then a coefficient is set to 0 and dropped

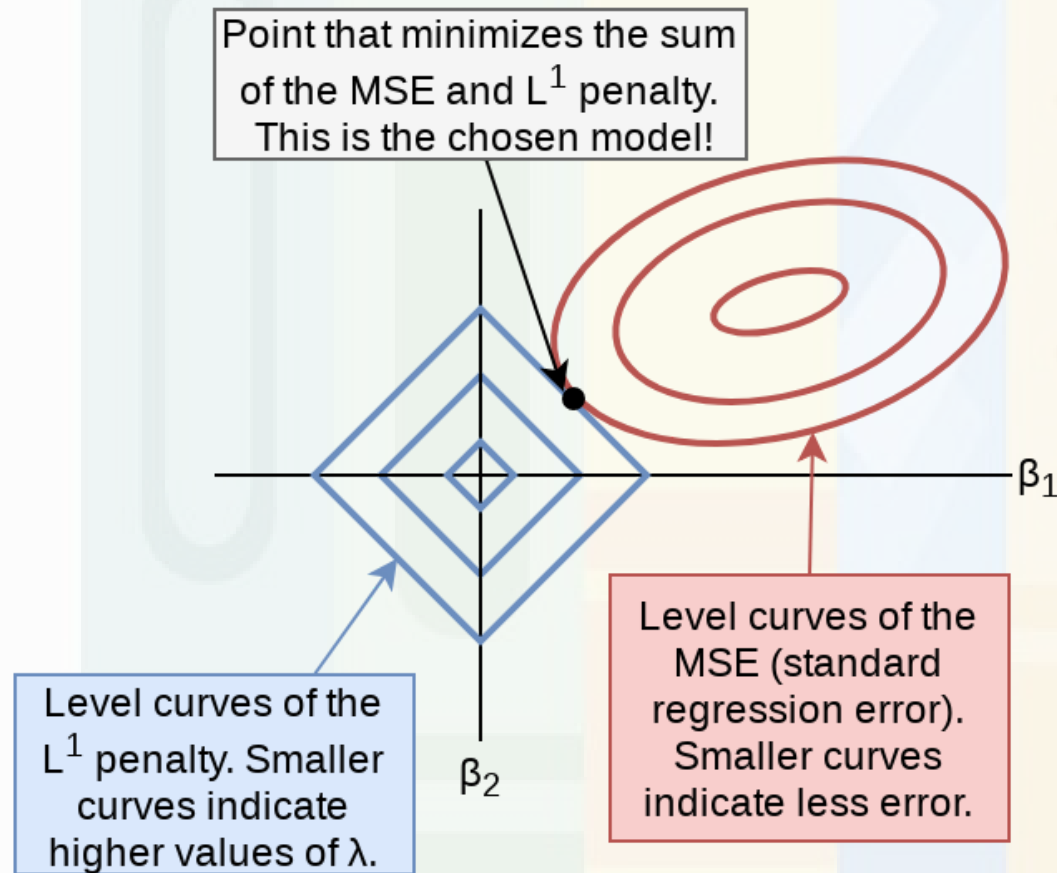
Illustration of LASSO in the *coefficient space* of a regression



What about other penalty types?

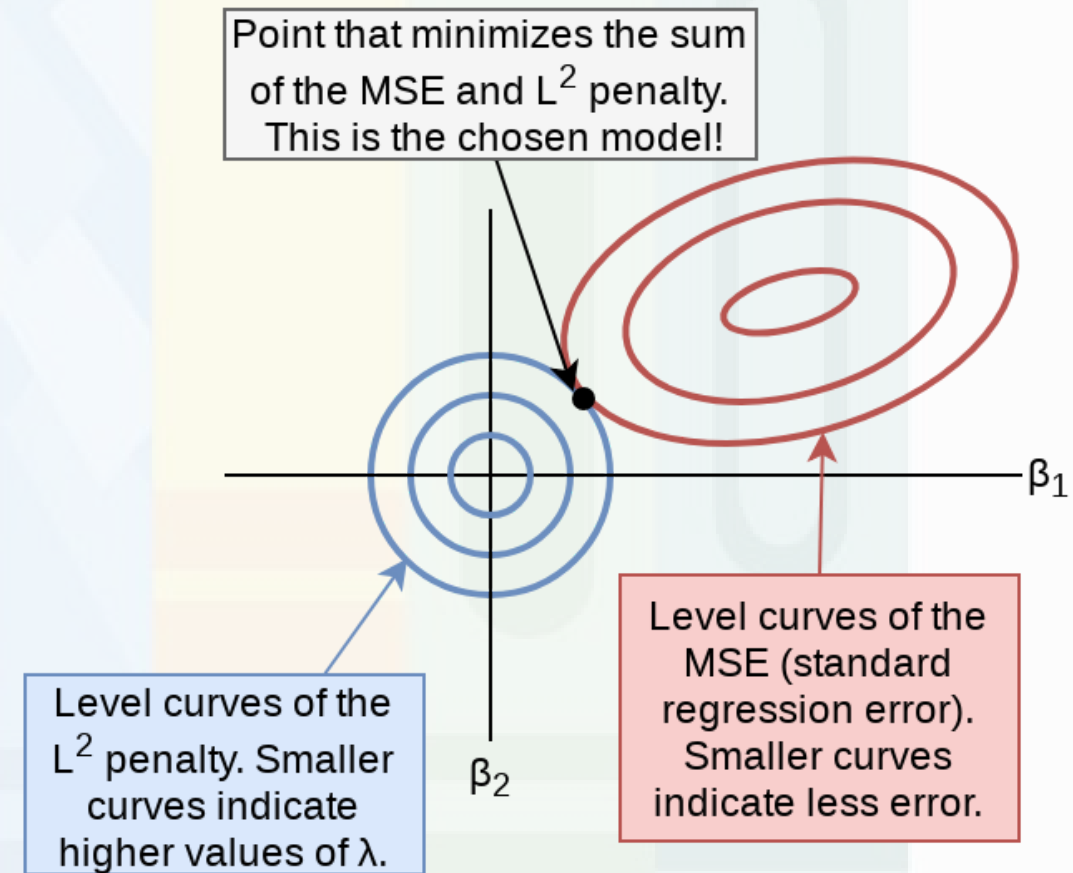
LASSO (L_1)

Illustration of LASSO in the *coefficient space* of a regression



Ridge (L_2)

Illustration of ridge in the *coefficient space* of a regression



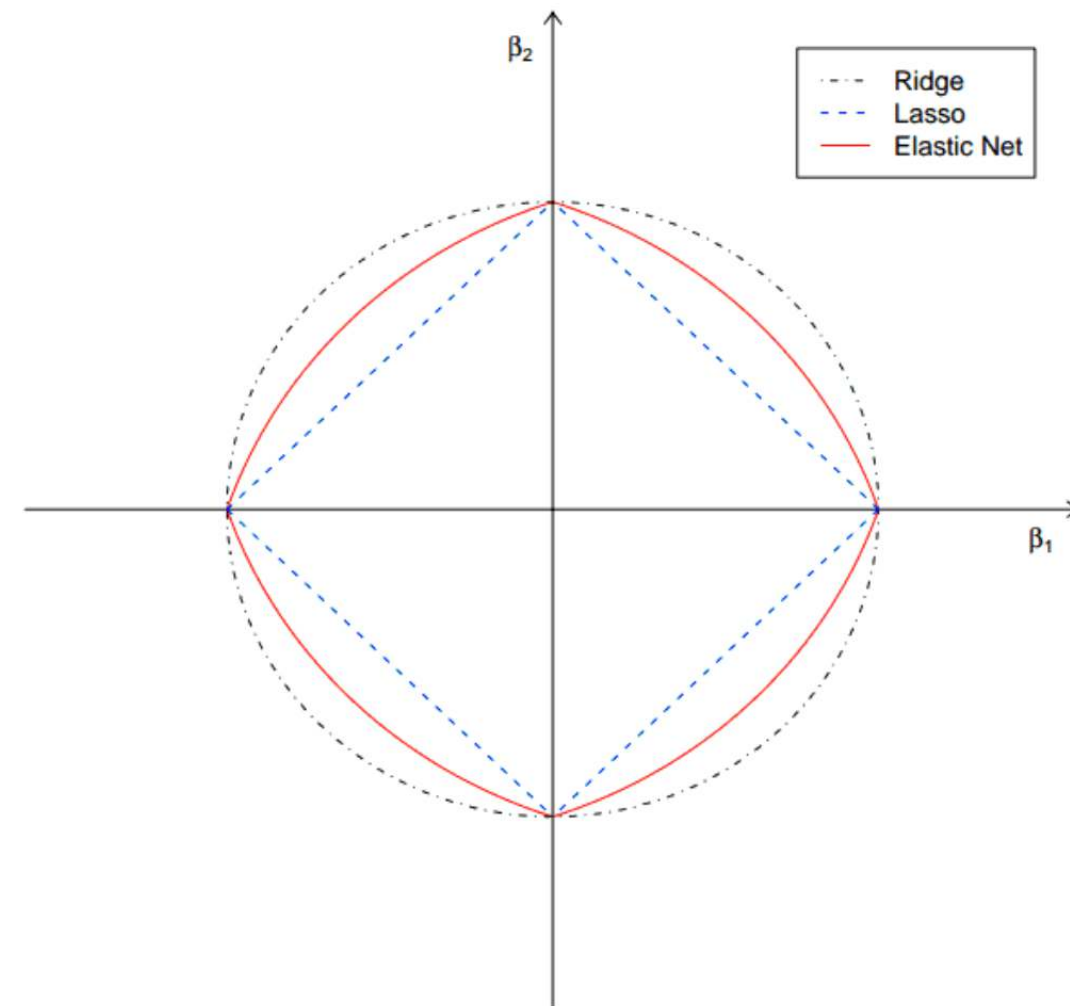
- Decreases coefficient values
 - Makes many of them 0
 - Increases prediction stability

- Decreases coefficient values
 - Increases prediction stability more
 - Less sensitive to outliers

Combining LASSO and Ridge: Elastic Net

- Elastic Net has both L_1 and L_2 penalties!
- Allows you to optimize the amount of selection effect you want from LASSO and the amount of shrinkage from Ridge
- A generalization of LASSO and Ridge

$$\min_{\beta \in \mathbb{R}} \left\{ \frac{1}{N} |\epsilon|_2^2 + \lambda_1 |\beta|_1 + \lambda_2 \|\beta\|^2 \right\}$$



Technical: Preparation



Importing data

- Python: We can use `pandas` to import the data set
- R: We can use `tidyverse` to import the data set
- Compressing a csv file can save 50-90% of the storage space of the file

```
df = pd.read_csv('../Data/S1_data.csv.gz')
```



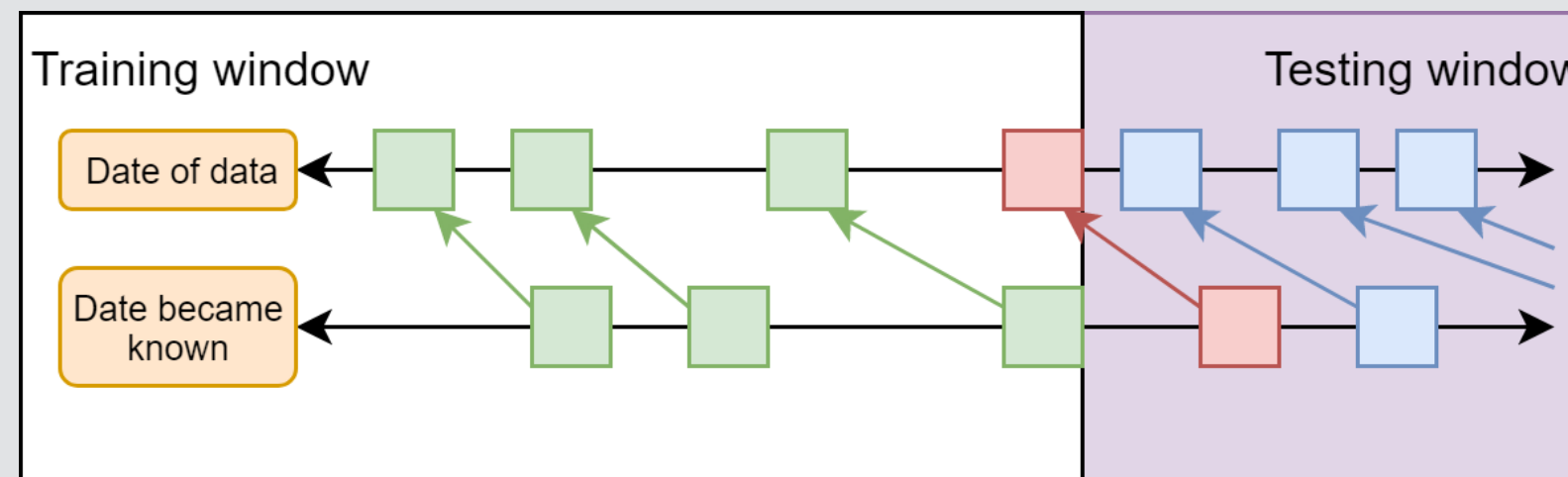
```
df = read_csv('../Data/S1_data.csv.gz')
```



- Note:
 - SAS, python pandas, and R can all handle `.csv.gz` and `.csv.zip` files
 - Stata is a bit tedious here, requiring uncompressing first
 - Either use your file manager or using Stata's `unzipfile` command

Validating predictive analyses

- Ideal:
 - Withhold the last year (or a few) of data when building the model
 - Check performance on *hold out sample*
 - This is *out of sample* testing
 - Ensure that the data is independent across time!



- Sometimes acceptable:
 - Withhold a random sample of data when building the model
 - Check performance on *hold out sample*
 - Potential problems with correlations between hold out sample and training sample

Training vs. testing split

- A simple approach is to split by time
- Check which years are in the data using `.unique()`

```
# Check the years in the data  
df['year'].unique()
```



```
## array([2002, 2003, 2004, 1999, 2000, 2001], dtype=int64)
```

```
# Check the years in the data  
unique(df$year)
```



```
## [1] 2002 2003 2004 1999 2000 2001
```

- Split out the last year as the testing sample
 - This can be done using a simple conditional
 - Final year is 2004, so...
 - Testing: `df.year == 2004`
 - Training: `df.year < 2004`

Splitting the sample

```
# Subset the final year to be the testing year  
train = df[df.year < 2004]  
test  = df[df.year == 2004]  
print(df.shape, train.shape, test.shape)
```

```
## (14301, 198) (11478, 198) (2823, 198)
```

```
# Subset the final year to be the testing year  
train <- df %>% filter(year < 2004)  
test  <- df %>% filter(year == 2004)  
print(c(nrow(df), nrow(train), nrow(test)))
```

```
## [1] 14301 11478 2823
```

- Note that the number of rows in `df` is the same as the sum of rows in `train` and `test`

Aside: Random testing sample

- In python, Scikit-learn (`sklearn`) can handle this robustly
 - Scikit-learn is a package focused on simple machine learning methods
 - Since random sampling is common in ML, Scikit-learn provides multiple ways to handle this.
 - The function is `sklearn.model_selection.train_test_split()`
 - Optionally you can stratify across classes in your data using the `stratify=` parameter
- In R, `caret` can handle this well using the `createDataPartition()` function

Technical: Running simple regressions

Package: Statsmodels

- The `statsmodels` package provides a suit of basic regression functions
- It supports most standard statistical approaches
 - OLS, Logit, GLM, Probit, Poisson, ARIMA, etc.
- It includes some other interesting functions as well, such as:
 - Imputation methods (e.g., MICE), GAMs, Quantile regression, Markov switching, etc.
- There are 2 interfaces to the package:
 1. `statsmodels.formula.api` (usually imported as `smf`) – pandas-friendly
 2. `statsmodels.api` (usually imported as `sm`) – requires data to be formatted differently

Linear regression (OLS)

- Unlike most statistical software, regressions in `statsmodels` require multiple steps.

Step 1: specify the regression structure

```
formula = 'sdvoll ~ ' + ' + '.join(vars_topic[0:-1])  
model = smf.ols(formula=formula, data=train)
```



- Note the use of `~` as the equals sign in the equation

Step 2: Run the regression

```
fit1 = model.fit()
```



Linear regression (OLS)

Step 3: Output the results (optional)

```
fit1.summary()
```



OLS Regression Results

Dep. Variable:	sdvol1	R-squared:	0.161
Model:	OLS	Adj. R-squared:	0.159
Method:	Least Squares	F-statistic:	73.45
Date:	Sun, 14 Aug 2022	Prob (F-statistic):	0.00
Time:	21:47:41	Log-Likelihood:	24508.
No. Observations:	11478	AIC:	-4.895e+04
Df Residuals:	11447	BIC:	-4.873e+04
Df Model:	30		
Covariance Type:	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
Intercept	0.0458	0.000	171.114	0.000	0.045	0.046
Topic_1_n_ol	1.1709	0.340	3.440	0.001	0.504	1.838

Base R

- Fitting regressions is straightforward in R

```
BD_eq <- as.formula(paste("sdvol1 ~ ", paste(paste0("Topic_", 1:30, "_n_oI"), collapse=" + "), collapse=""))  
model <- lm(BD_eq, train)  
summary(model)
```



```
##  
## Call:  
## lm(formula = BD_eq, data = train)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max   
## -0.18799 -0.01707 -0.00646  0.00904  0.49410   
##  
## Coefficients:  
##              Estimate Std. Error t value Pr(>|t|)      
## (Intercept)   0.0457521  0.0002674 171.114 < 2e-16 ***  
## Topic_1_n_oI  1.1709484  0.3404372   3.440 0.000585 ***  
## Topic_2_n_oI  0.5367261  0.2615383   2.052 0.040174 *  
## Topic_3_n_oI  0.4004462  0.4160324   0.963 0.335801   
## Topic_4_n_oI  0.6475066  0.2386256   2.713 0.006668 **  
## Topic_5_n_oI  0.6776698  0.2462900   2.752 0.005941 **  
## Topic_6_n_oI  0.5421747  0.3630189   1.494 0.135330   
## Topic_7_n_oI -0.6519468  0.2858123  -2.281 0.022565 *  
## Topic_8_n_oI  0.5089414  0.2529234   2.012 0.044219 *  
## Topic_9_n_oI  2.1940373  0.2245302   9.772 < 2e-16 ***  
## Topic_10_n_oI 0.6721560  0.2073181   3.242 0.001190 **
```


Logistic regression in Python

```
formula = 'Restate_Int ~ ' + \  
          ' + '.join(vars_financial) + ' + ' +\  
          ' + '.join(vars_style) + ' + ' +\  
          ' + '.join(vars_topic[0:-1]) # Drop the final value to avoid multicollinearity  
model = smf.logit(formula=formula, data=train)  
fit_logit = model.fit()
```

```
## Warning: Maximum number of iterations has been exceeded.  
##           Current function value: 0.054196  
##           Iterations: 35  
##  
## M:\Python_environments\Teaching_ML_v1\lib\site-packages\statsmodels\discrete\discrete_model.py:1810: RuntimeWarning: overfl  
##   return 1/(1+np.exp(-X))  
## M:\Python_environments\Teaching_ML_v1\lib\site-packages\statsmodels\base\model.py:568: ConvergenceWarning: Maximum Likelihc  
##   ConvergenceWarning)
```

```
fit_logit.summary()
```

Logit Regression Results

Dep. Variable:	Restate_Int	No. Observations:	11478
Model:	Logit	Df Residuals:	11410
Method:	MLE	Df Model:	67
Date:	Sun, 14 Aug 2022	Pseudo R-squ.:	0.1205
i :	. .	- i li :	-

Logistic regression in R

```
BCE_eq <- as.formula(paste("Restate_Int ~ logtotasset + rsst_acc + chg_recv + chg_inv +  
soft_assets + pct_chg_cashsales + chg_roa + issuance +  
oplease_dum + book_mkt + lag_sdvol + merger + bigNaudit +  
midNaudit + cffin + exfin + restruct + bullets + headerlen +  
newlines + alltags + processedsize + sentlen_u + wordlen_s +  
paralen_s + repetitious_p + sentlen_s + typetoken +  
clindex + fog + active_p + passive_p + lm_negative_p +  
lm_positive_p + allcaps + exclamationpoints + questionmarks + ",  
paste(paste0("Topic_",1:30,"_n_oI"), collapse=" + "), collapse=""))  
  
model_logit <- glm(BCE_eq, train, family="binomial")
```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

```
summary(model_logit)
```

```
##  
## Call:  
## glm(formula = BCE_eq, family = "binomial", data = train)  
##  
## Deviance Residuals:  
##      Min       1Q   Median       3Q      Max   
## -1.1728  -0.1591  -0.1092  -0.0739   3.6910   
##  
## Coefficients:  
##              Estimate Std. Error z value Pr(>|z|)   
## (Intercept)  -6.634e+00  5.591e+00  -1.187  0.23541   
## logtotasset   9.363e-02  6.442e-02   1.454  0.14607   
## rsst_acc      3.269e-01  3.226e-01   1.013  0.31095   
## chg_recv      6.838e-01  1.307e+00   0.523  0.60085   
## chg_inv     -1.428e+00  1.509e+00  -0.947  0.34378   
## soft_assets   1.451e+00  4.698e-01   3.088  0.00201 **
```

Technical: Measuring predictive performance



Linear predictive power

- 2 methods that are often used are:
 - RMSE: Root Mean Squared Error
 - MAE: Mean Absolute Error

RMSE

```
sklearn.metrics.mean_squared_error()
```

```
apply_rmse <- function(v1, v2) {  
  sqrt(mean((v1 - v2)^2, na.rm=T))  
}
```



MAE

```
sklearn.metrics.mean_absolute_error()
```

```
apply_mae <- function(v1, v2) {  
  mean(abs(v1-v2), na.rm=T)  
}
```

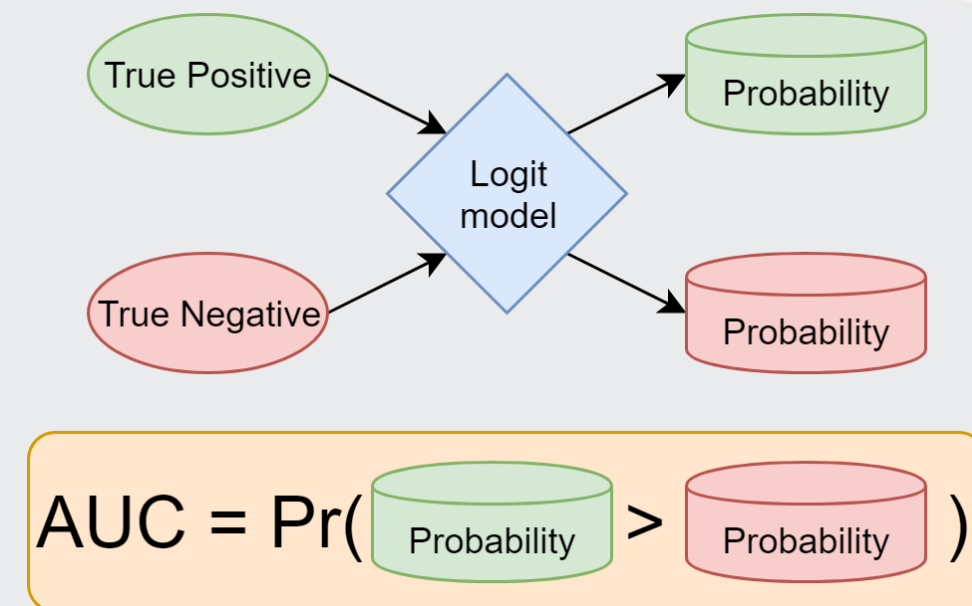


Logistic predictive power

- For logistic regression, ROC AUC is a good measure
- Use `sklearn` in python or `yardstick` in R

```
Y_hat_test = fit_logit.predict(test)  
auc = metrics.roc_auc_score(test.Restate_Int, Y_hat_test)
```

```
test$Y_hat_test <- predict(model_logit, test, type="response")  
auc_out <- test %>% roc_auc(truth=as.factor(Restate_Int),
```

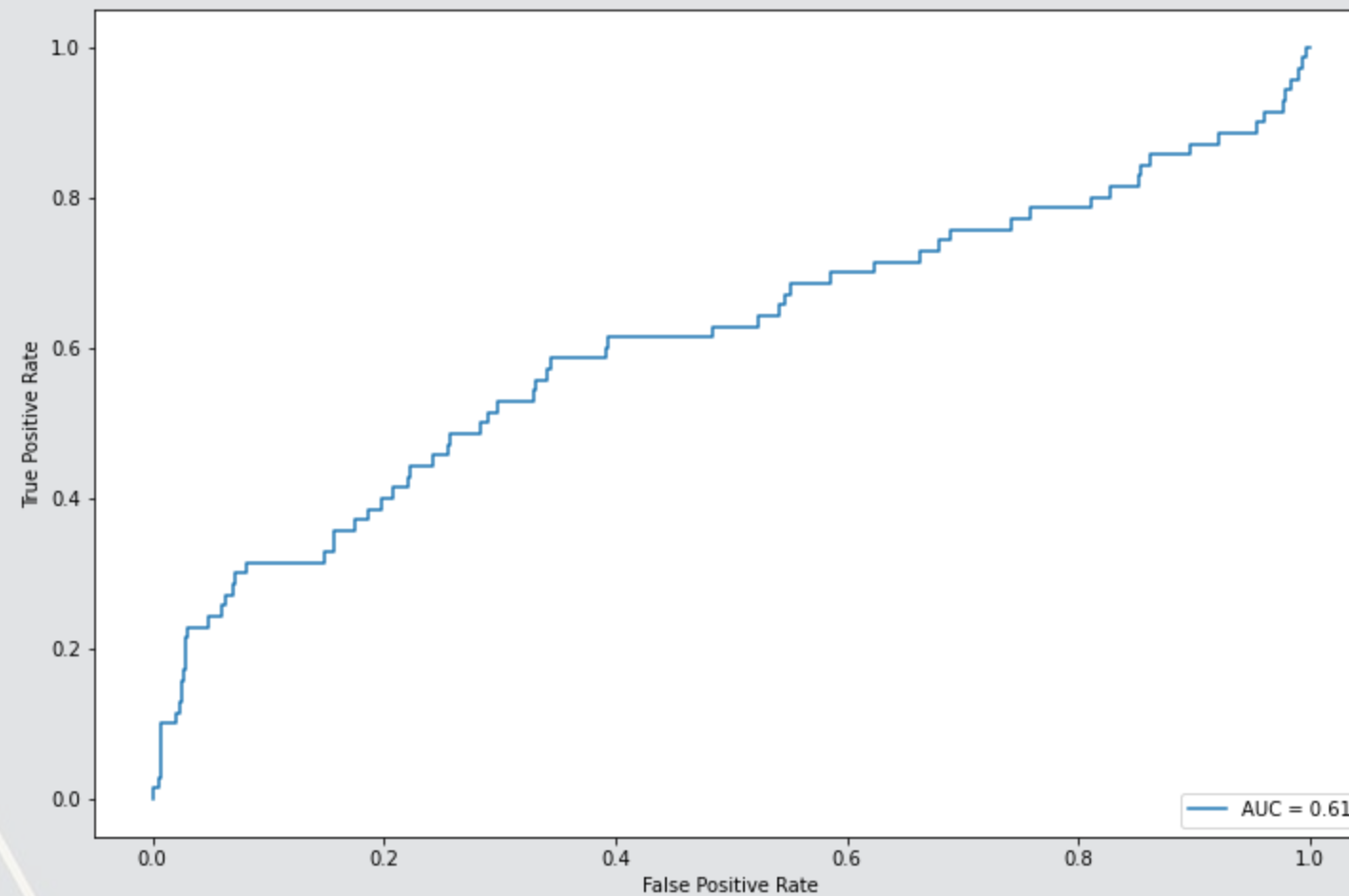


Visualizing AUC with the ROC curve

- `sklearn` makes it easy to output a ROC curve as well

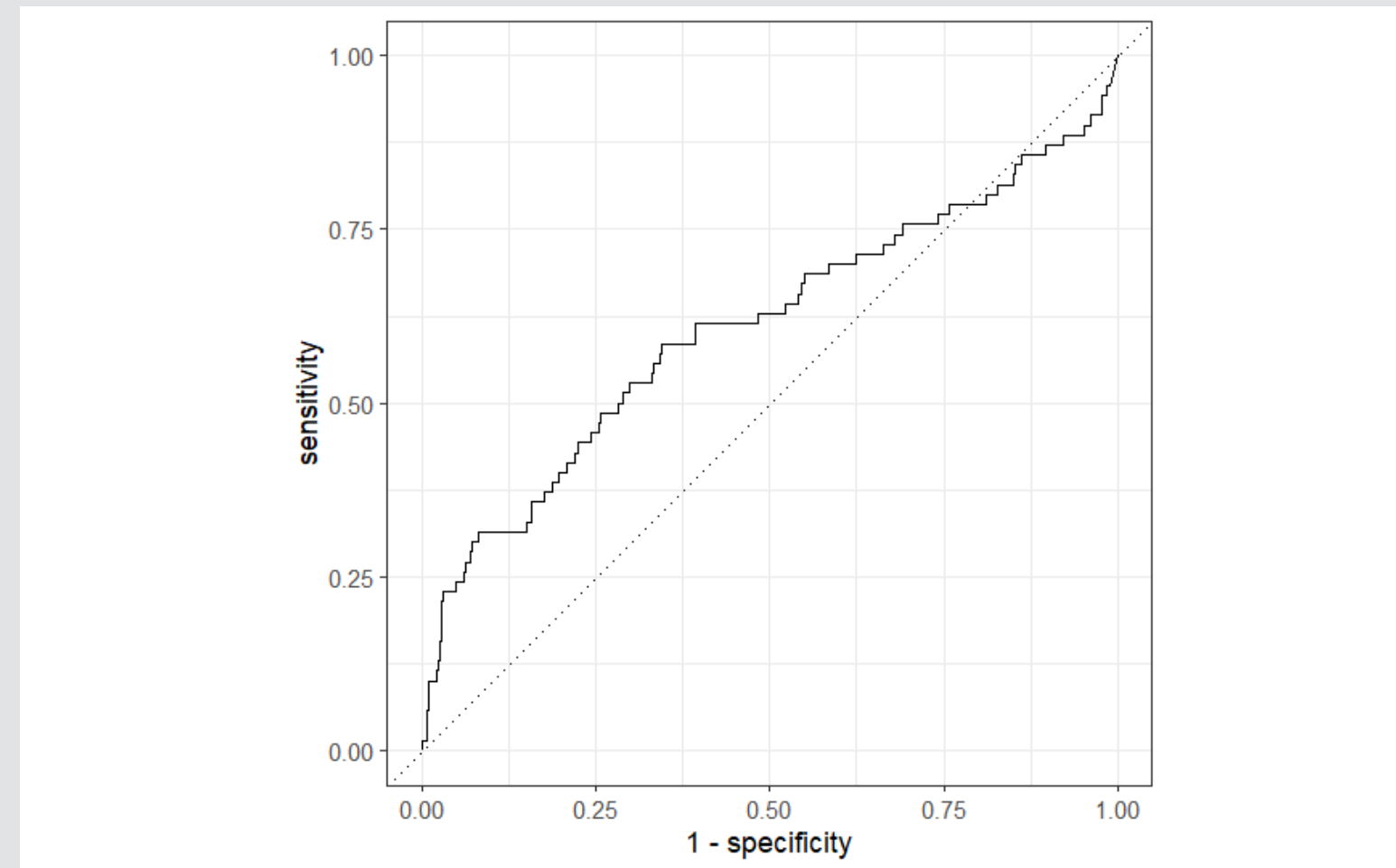
```
# Logit, out-of-sample
Y_hat_test = fit_logit.predict(test)
auc = metrics.roc_auc_score(test.Restate_Int, Y_hat_test)

fpr, tpr, thresholds = metrics.roc_curve(test.Restate_Int, Y_hat_test)
display = metrics.RocCurveDisplay(fpr=fpr, tpr=tpr, roc_auc=auc)
display.plot()
```



Using yardstick in R

```
test$Y_hat_test <- predict(model_logit, test, type="response")
test %>%
  roc_curve(truth=as.factor(Restate_Int), estimate=Y_hat_test, event_level='second') %>%
  autoplot()
```



Technical: Implementing LASSO

Using python: Setting up to use Scikit-Learn

- Scikit-learn, like many machine learning packages, expects separate data sets or matrices for DVs and IVs
- LASSO, Ridge, and Elastic net are also particular about data format:

Every input should be normalized to a Z-score! (python-specific requirement)

- Scikit-learn has this all built in, so it will be easy

```
vars = vars_topic
scaler_X = preprocessing.StandardScaler()
scaler_X.fit(train[vars])
train_X_linear = scaler_X.transform(train[vars])
test_X_linear = scaler_X.transform(test[vars])
```



- `sklearn.preprocessing.StandardScaler()` defaults to transforming to Z-scores
- Applying `.fit()` with data makes it calculate the mean and standard deviation of each column
- Applying `.transform()` with data applies the Z-score based on the fitted parameters
 - Avoids any look-ahead bias in our testing sample!

Using Python: Setting up to use Scikit-Learn

```
scaler_Y = preprocessing.StandardScaler()
scaler_Y.fit(np.array(train.sdvoll).reshape(-1, 1))
train_Y_linear = scaler_Y.transform(np.array(train.sdvoll).reshape(-1, 1))
test_Y_linear = scaler_Y.transform(np.array(test.sdvoll).reshape(-1, 1))
```



- Inputs are required to be 2D matrices by `sklearn`
- The `np.array(____).reshape(-1, 1)` bit is to cast the Pandas series back into a 2D matrix -
`np.array()` casts the pandas series object to an array (matrix), but it is only 1D
 - `.reshape(-1, 1)` forces the matrix to be a column (and thus 2D) instead of a 1D row matrix

Using Python: Simple LASSO, linear

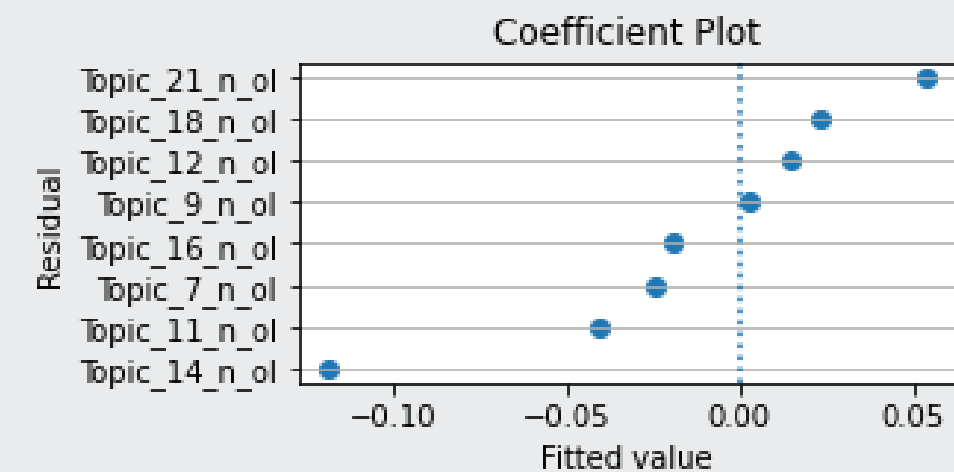
Fitting a LASSO with a pre-specified penalty

```
reg_lasso = linear_model.Lasso(alpha=0.1)  
reg_lasso.fit(train_X_linear, train_Y_linear)
```



Custom coefficient plot function

```
coefplot(vars, reg_lasso.coef_)
```



Not too difficult, but the coefplot function is custom (see Jupyter notebook for it)

Using R: Setting up to use glmnet

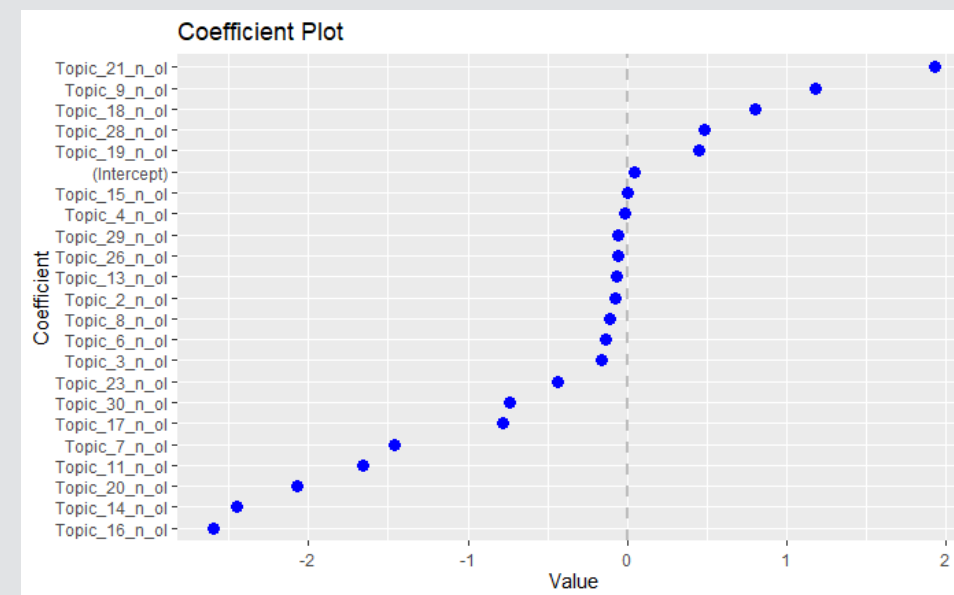
- The `glmnet` expects data as separate matrices for X and Y measures
- It does not require data to be Z-scores – it is invariant to this
- The `model.matrix()` and `model.frame()` commands from Base R make this easy

```
x_lm <- model.matrix(BD_eq, data=train)[,-1] # [-1] to remove intercept  
y_lm <- model.frame(BD_eq, data=train)[,"sdvoll"]
```



Using R: Running glmnet

```
fit_LASSO_lm <- glmnet(x=x_lm, y=y_lm,  
                      family = "gaussian",  
                      alpha = 1 # Specifies LASSO. alpha = 0 is ridge  
                      )  
  
coefplot(fit_LASSO_lm, sort='magnitude')
```



In this case, `coefplot` is available from CRAN

Cross validation

What is cross validation?

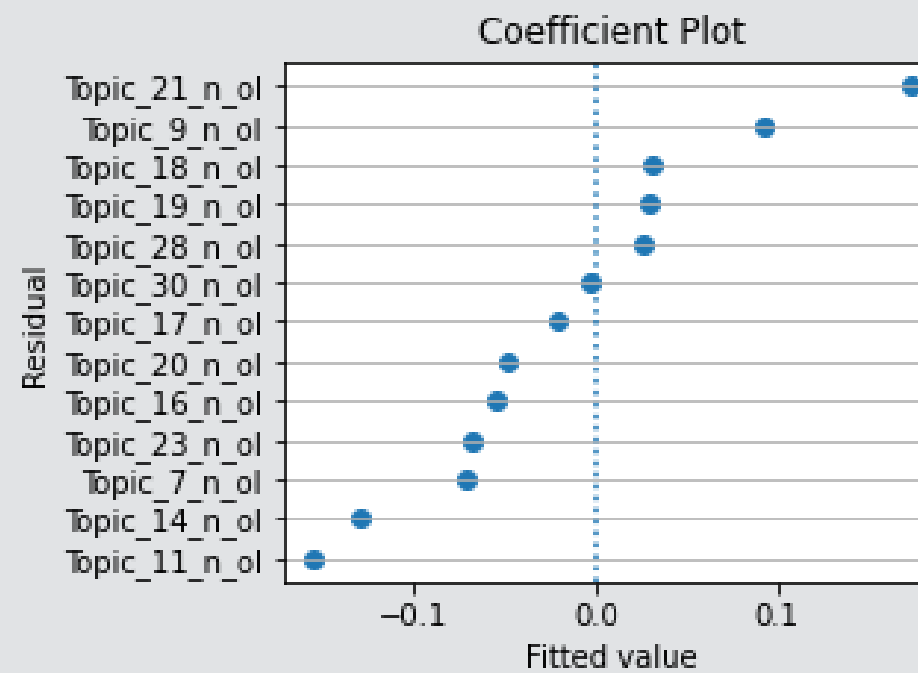
- Validation is where you keep part of the training sample as a hold out sample to evaluate and improve your algorithm against
 - This prevents biasing towards the real hold out sample (the testing sample)
- Cross validation takes this further by making a bunch of validation samples,
- An example of 10-fold cross validation:
 1. Randomly splits the data into 10 groups
 2. Runs the algorithm on 90% of the data ($10 - 1 = 9$ groups)
 3. Determines the best model based on the performance of the group that was left out
 4. Repeat steps 2 and 3 $10 - 1 = 9$ more times
 5. Uses the best overall model across all 10 hold out samples

Scikit-learn has this built in! So does glmnet!

10-fold CV LASSO, linear, Python

```
reg_lasso = linear_model.LassoCV(cv=10)
reg_lasso.fit(train_X_linear, np.ravel(train_Y_linear))

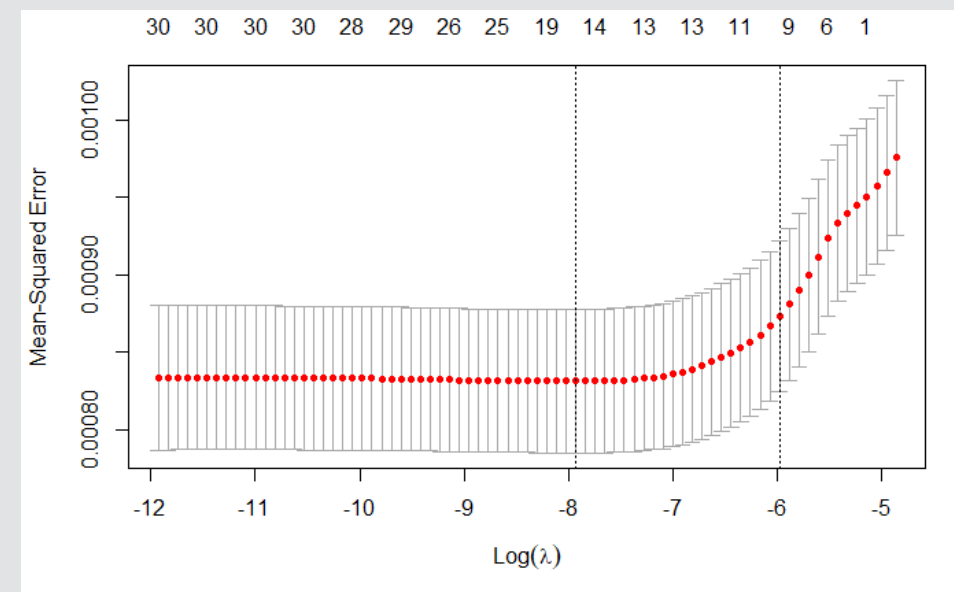
coefplot(vars, reg_lasso.coef_)
```



10-fold CV LASSO, linear, R

- To replicate our linear LASSO:

```
cvfit_lm = cv.glmnet(x=x_lm, y=y_lm, family = "gaussian", alpha = 1, type.measure="mse")  
plot(cvfit_lm)
```



Note: This is optimizing MSE instead of R^2 – glmnet doesn't support R^2 !

10-fold CV elastic net, linear, Python

- Need to specify values to examine for the ratio between L_1 and L_2 penalty
 - `l1_ratio=1` is a LASSO, `l1_ratio=0` is Ridge, in between is elastic net

```
reg_EN = linear_model.ElasticNetCV(cv=10, l1_ratio=[.1, .5, .7, .9, .95, .99, 1])  
reg_EN.fit(train_X_linear, np.ravel(train_Y_linear))
```



Note: This does CV over both parameters!

10-fold CV elastic net, linear, R

- In R, `glmnet` can do this too
 - `lambda=1` is LASSO
 - `lambda=0` is Ridge
 - If `lambda` is set between 0 and 1, it's an elastic net!
- To replicate our linear LASSO:

```
cvfit_en = cv.glmnet(x=x, y=y, family = "binomial", alpha = 0.5, type.measure="auc")
```



Note: This does CV only over the penalty parameter. You need to build your own grid over the alpha parameter

Conclusion



Wrap-up

Econometrics

- R and Stata are both better for this, python is capable but not as simple

Machine learning regression

- Python is better at this than basic regression
- In some circumstances, these techniques are
 - More econometrically defensible
 - More robust
 - More accurate
- R's `glmnet` package is more efficient and easier to use
 - But the elastic net implementation is more flexible for CV in Python
- Stata has an interesting implementation in `lassopack`
- For more interesting variants, check out R's `hdm`

Packages used for these slides

Python

- matplotlib
- numpy
- pandas
- scikit-learn
- statsmodels

R

- glmnet
- kableExtra
- knitr
- reticulate
- revealjs
- tidyverse
- yardstick

References

- Bao, Yang, and Anindya Datta. “Simultaneously discovering and quantifying risk types from textual risk disclosures.” *Management Science* 60, no. 6 (2014): 1371-1391.
- Brown, Nerissa C., Richard M. Crowley, and W. Brooke Elliott. “What are you saying? Using topic to detect financial misreporting.” *Journal of Accounting Research* 58, no. 1 (2020): 237-291.
- Chahuneau, Victor, Kevin Gimpel, Bryan R. Routledge, Lily Scherlis, and Noah A. Smith. “Word salad: Relating food prices and descriptions.” In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pp. 1357-1367. 2012.
- Mullainathan, Sendhil, and Jann Spiess. “Machine learning: an applied econometric approach.” *Journal of Economic Perspectives* 31, no. 2 (2017): 87-106.

Custom code

```
# Replication of R's coefplot function for use with sklearn's linear and logistic LASSO

def coefplot(names, coef, title=None):
    # Make sure coef is list, cast to list if needed.
    if isinstance(coef, np.ndarray):
        if len(coef.shape) > 1:
            coef = list(coef[0])
        else:
            coef = list(coef)

    # Drop unneeded vars
    data = []
    for i in range(0, len(coef)):
        if coef[i] != 0:
            data.append([names[i], coef[i]])
    data.sort(key=lambda x: x[1])

    # Add in a key for the plot axis
    data = [data[i] + [i+1] for i in range(0, len(data))]

    fig, ax = plt.subplots(figsize=(4, 0.25*len(data)))

    ax.scatter([i[1] for i in data], [i[2] for i in data])

    ax.grid(axis='y')
    ax.set(xlabel="Fitted value", ylabel="Residual", title=(title if title is not None else "Coefficient Plot"))

    ax.axvline(x=0, linestyle='dotted')
    ax.set_yticks([i[2] for i in data])
    ax.set_yticklabels([i[0] for i in data])

    return ax
```

